

AD-A068 544

NAVAL POSTGRADUATE SCHOOL MONTEREY CALIF

F/G 12/1

DEVELOPMENT OF CLUSTER ANALYSIS METHODS SUITABLE FOR STUDENT OP--ETC(U)

MAR 79 J W AIKEN

UNCLASSIFIED

NL

1 OF 2

AD
A068544



② LEVEL II

NAVAL POSTGRADUATE SCHOOL
Monterey, California

AD A068544



DDC
RECEIVED
MAY 15 1979
B

DDC FILE COPY

THESIS

DEVELOPMENT OF CLUSTER ANALYSIS METHODS
SUITABLE FOR STUDENT OPINION DATA

by

Joel Weston Aiken

March 1979

Thesis Advisor:

R.R. Read

Approved for public release; distribution unlimited.

70 05 11 064

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM
1. REPORT NUMBER	2. GOVT ACCESSION NO.	3. RECIPIENT'S CATALOG NUMBER
4. TITLE (and Subtitle)		5. TYPE OF REPORT & PERIOD COVERED
Development of Cluster Analysis Methods Suitable for Student Opinion Data		Master's Thesis, March 1979
7. AUTHOR(s)		6. PERFORMING ORG. REPORT NUMBER
Joel Weston/Aiken		8. CONTRACT OR GRANT NUMBER(s)
9. PERFORMING ORGANIZATION NAME AND ADDRESS		10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS
Naval Postgraduate School Monterey, California 93940		
11. CONTROLLING OFFICE NAME AND ADDRESS		12. REPORT DATE
Naval Postgraduate School Monterey, California 93940		March 1979
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office)		13. NUMBER OF PAGES
		151
		15. SECURITY CLASS. (of this report)
		Unclassified
		15a. DECLASSIFICATION/DOWNGRADING SCHEDULE
16. DISTRIBUTION STATEMENT (of this Report)		
Approved for public release; distribution unlimited.		
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)		
18. SUPPLEMENTARY NOTES		
19. KEY WORDS (Continue on reverse side if necessary and identify by block number)		
Cluster Analysis		
20. ABSTRACT (Continue on reverse side if necessary and identify by block number)		
<p>The Naval Postgraduate School's Student Opinion Form data were subjected to study through the use of two cluster analysis techniques: (1) K-MEANS partitioning method and (2) Chernoff's FACES. Much developmental work was performed to tailor these methods to the special requirements of the data set. A thorough multivariate statistical review provided the basis for choosing optimality criteria and distance functions for use in the MIKCA</p>		

DD FORM 1473
1 JAN 73
(Page 1)EDITION OF 1 NOV 65 IS OBSOLETE
S/N 0102-014-6601

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

251 450

1

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE(When Data Entered)

#20 (ABSTRACT) - Cont'd

4 (Multivariate Iterative K-MEANS Clustering Algorithm). Alterations were made to the computer code to allow the analysis to include the effect of class size on cluster membership. Use of the linear discriminant function aided in identifying variables for use in constructing features of the computer-drawn faces. This approach to the Chernoff's FACES technique shows promise but needs further development. A principal components analysis of the data showed it to be essentially one dimensional. Partitioning the data into four clusters shows that the scoring of the courses varies inversely with class size.

ACCESSION for	
NTIS	White Section <input checked="" type="checkbox"/>
DDC	Buff Section <input type="checkbox"/>
UNANNOUNCED	<input type="checkbox"/>
JUSTIFICATION	
BY	
DISTRIBUTION/AVAILABILITY CODES	
Dist.	APPROPRIATE or SPECIAL
A	

DD Form 1473
1 Jan 73
S/N 0102-014-6601

UNCLASSIFIED

2 SECURITY CLASSIFICATION OF THIS PAGE(When Data Entered)

Approved for public release; distribution unlimited.

Development of Cluster Analysis Methods
Suitable for Student Opinion Data

by

Joel Weston Aiken
Lieutenant Commander, United States Navy
B.S., University of North Carolina at Chapel Hill, 1969

Submitted in partial fulfillment of the
requirements for the degree of

MASTER OF SCIENCE IN OPERATIONS RESEARCH

from the

NAVAL POSTGRADUATE SCHOOL

March 1979

Author

Approved by:

R R Reed

Thesis Advisor

John R. Bonney

Second Reader

Michael D. Foreman
Chairman, Department of Operations Research

A. Shady
Dean of Information and Policy Sciences

ABSTRACT

The Naval Postgraduate School's Student Opinion Form data were subjected to study through the use of two cluster analysis techniques: (1) K-MEANS partitioning method and (2) Chernoff's FACES. Much developmental work was performed to tailor these methods to the special requirements of the data set. A thorough multivariate statistical review provided the basis for choosing optimality criteria and distance functions for use in the MIKCA (Multivariate Iterative K-MEANS Clustering Algorithm). Alterations were made to the computer code to allow the analysis to include the effect of class size on cluster membership. Use of the linear discriminant function aided in identifying variables for use in constructing features of the computer-drawn faces. This approach to the Chernoff's FACES technique shows promise but needs further development. A principal components analysis of the data showed it to be essentially one dimensional. Partitioning the data into four clusters shows that the scoring of the courses varies inversely with class size.

TABLE OF CONTENTS

I.	INTRODUCTION -----	10
II.	CLUSTER ANALYSIS -----	14
	A. ORIGIN AND THEORY -----	14
	B. SCATTER MATRIX DECOMPOSITION -----	18
	C. OPTIMALITY CRITERIA -----	20
	D. DISTANCE CONSIDERATIONS -----	24
III.	THE DATA SET -----	32
	A. ORIGIN -----	32
	B. TRANSFORMATION -----	35
	C. PRINCIPAL COMPONENTS -----	37
IV.	THE MIKCA METHOD -----	44
	A. THE ALGORITHM -----	44
	B. MODIFIED MIKCA -----	48
	C. RESULTS -----	52
V.	DISCRIMINANT ANALYSIS -----	59
	A. THEORY -----	59
	B. RESULTS -----	60
VI.	CHERNOFF FACES -----	66
	A. BACKGROUND -----	66
	B. FEATURE-VARIABLE RELATIONSHIP -----	72
	C. CLUSTERING THE FACES -----	76
	D. PROBLEMS ENCOUNTERED -----	82
VII.	COMPARISON COEFFICIENTS -----	85
	A. ALTERNATIVES -----	85
	B. THE TECHNIQUE -----	87

VIII. SUMMARY AND CONCLUSIONS -----	89
APPENDIX A: TEST STATISTIC FOR HOMOGENEITY OF DISPERSIONS -----	94
APPENDIX B: PROFILE CHARTS -----	96
APPENDIX C: FACES CONSTRUCTION FORMULAE -----	103
APPENDIX D: EE DATA FROM QUARTER 781 -----	109
APPENDIX E: EXAMPLE OF COMPARISON COEFFICIENT -----	111
APPENDIX F: MODIFIED MIKCA COMPUTER PROGRAM COMPARISON COEFFICIENT COMPUTER PROGRAM -----	114
BIBLIOGRAPHY -----	149
INITIAL DISTRIBUTION LIST -----	151

LIST OF TABLES

1. DATA CATEGORIES -----	34
2. PRINCIPAL COMPONENTS ANALYSIS -----	41
3. COMPARISON COEFFICIENTS FOR VARIOUS SOLUTIONS -----	54
4. RESULTS OF DISCRIMINANT ANALYSIS -----	62
5. RANGES AND DESCRIPTION OF FEATURES -----	71
6. FEATURE-VARIABLE COMBINATIONS -----	74
7. COMPARISON COEFFICIENTS FOR JUDGES -----	79
8. CONTINGENCY TABLE -----	85
9. COMPARISON COEFFICIENT EXAMPLE -----	87

LIST OF FIGURES

1. GEOMETRIC DISTANCE CONSIDERATIONS -----	24
2. SOF FORM -----	33
3. TRANSFORMATIONS -----	38
4. WITHIN COURSE VARIABILITY VS. TIME ON BOARD -----	43
5. MIKCA FLOW CHART -----	45
6. TRACE W VS. g -----	57
7. DET W VS. g -----	58
8. DISCRIMINANT SPACE -----	63
9. BRUCKNER'S OIL COMPANY FACES -----	69
10. FACE CONSTRUCTION -----	70
11. CLUSTERS OF 33 FACES FROM EE DEPT, QTR 781 -----	77
12. EXPERIMENTAL FACES -----	81

ACKNOWLEDGMENTS

The author would like to acknowledge the assistance of Mrs. Pat Meadows who organized the data from card decks into an accessible data disc so that retrieval was made easy and the privacy of individual instructors was respected. Professor G. Lindsay's work in scale analysis on part of the data was used to help select transformations to stabilize the variability. Professor J. Hartman suggested the sum of squares technique which led to the development of the comparison measure. Professor D. Kirk (Chairman of the EE Dept) participated in clustering the faces and gave a detailed description of how he formed the clusters.

The following students participated as judges and offered comments on their clustering methods: LCDR John Scott Redd, LCDR Raymond J. Morris, LT James Kevin McDermott and LT Howard S. Hilley. Mrs. Joel W. Aiken also acted as a judge, assisted with typing the first draft, and provided moral support. Professor R.R. Read gave willingly of his time in providing a background in multivariate statistics and guidance in the research effort.

I. INTRODUCTION

The Student Opinion Form (SOF) used at the Naval Postgraduate School provides an organized information gathering mechanism about each course (and its instructor) as perceived by the students. The information obtained from the SOF data is used for administrative review of faculty performance and for feedback to the instructor to aid in self-development. The former use is hampered by the fact that the data are multivariate in nature and represent a complicated set of interactions between the instructor's performance, the nature of the course, and the group of students. There is need for methodology which can disentangle those interactions and summarize the data in a meaningful way.

It is the purpose of this thesis to develop suitable cluster analysis methods for studying the data and discovering any hidden structure they may possess. Concurrently, a certain amount of exploratory data analysis took place, and those results are reported also.

At the completion of every quarter, students are requested to respond to a 16-item SOF questionnaire for each course in which they are enrolled. The data are viewed as an n by p matrix, representing n observations (SOFs), each of which is measured on p (16) different variables. For this research the mean vector of each course was computed. Then attempts

were made to discover natural clusters of these mean vectors which in turn can be interpreted as the underlying structure in the data. Since the number of students per course is quite variable, the mean vectors are not equally well determined. Also, the matrix of mean vectors may have a covariance structure quite different from that of the full n by p data matrix.

The clustering objective was pursued by two multivariate statistical methods: one computer-graphic technique referred to as Chernoff's FACES, and a second, more mathematically oriented approach called K-MEANS. The former produces computer-drawn cartoon faces, the features of which are controlled by variables in the data. The assignment of variables to features was aided by the use of linear discriminant analysis. One face is produced from each course mean vector, and then the researcher is able to study the appearance of the faces and cluster together those that display similar characteristics. The second method utilizes a computer program called MIKCA (Multivariate Iterative K-MEANS Clustering Algorithm) which is based on the K-MEANS method. It forms an initial partition of the data and then transfers observations between clusters in order to improve an optimality criterion function. In this iterative manner, MIKCA ultimately stabilizes and provides an "optimal" cluster solution.

In addition a modified MIKCA technique was employed. Alterations were made to the basic computer code to enable

the program to weight each mean vector by the number of students in the course. This modification may be likened to a one-way Analysis of Variance (ANOVA) having unbalance in the number of observations per treatment. The result is to stabilize the relative variability of the various course mean vectors.

Most multivariate analysis methodology is derived assuming the data have a multivariate normal distribution with common covariance matrix. The performance of the MIKCA program and the linear discriminant analysis will not depend greatly upon this assumption provided the clusters are well defined. On the other hand, if the clusters are not well separated, the results of the programs will be sensitive to these assumptions, and this is the condition anticipated. Accordingly, a transformation was sought toward this end. The one selected is essentially a logistic function.

It is frequently necessary to compare the agreement of cluster solutions produced under different conditions or by different methods. For this purpose, a computer program was written which provides an ad hoc measure of the amount of agreement between the results of two or more solutions. A number between zero and one, called the comparison coefficient, is the resulting measure of association.

This thesis was largely exploratory and should serve as a firm foundation for future study of the SOF data in particular and similarly structured multivariate data sets

in general. A number of unexpected questions are raised during the exploratory phases of this research. It was not possible to answer many of these questions, and their consideration is left to other researchers. During the development of the methodologies, some new and challenging problems were encountered. Many of these had to be given rather short treatment in the interest of meeting the original objectives. It should be emphasized that although some very interesting facts are revealed in this thesis, the results are by no means considered to describe completely the information hidden in the data.

II. CLUSTER ANALYSIS

A. ORIGIN AND THEORY

Cluster analysis is the name given to a body of diverse techniques for discovering taxonomical structure within bodies of data. It is one of several methodologies included in the broader category called classification. In cluster analysis little or nothing is known about the category structure. All that is available is a collection of observations whose category memberships are unknown, and one must discover a category structure which fits the observations. The objective is to find the natural groups by sorting the observations such that the association is high among members of the same group and low between members of different groups. The great challenge to the researcher is finding the most appropriate way of defining "natural groups" and "association." Cluster analysis is closely related to and often confused with discriminant analysis, a statistical procedure for assigning new observations to known groups. In contrast to discriminant analysis, clustering refers to discovery of the initial groups.

Although modern clustering techniques began development in biological taxonomy, they are generally applicable to all types of data. Any method which partitions a set of objects into subsets on the basis of measurements taken on every object qualifies as a clustering method. Cluster

analysis techniques are most often applied in multivariate settings, that is where each of n observations is measured on p different variables. A clear intuitive picture of the concept is helpful in appreciating the value of cluster analysis and the situations to which it might be applied. In a geometric sense, every object (observation) may be viewed as a point in p -dimensional Euclidean space. This swarm of data points may contain dense regions or "clouds" of data points which are separable from other regions containing a low density of points. These denser regions constitute what are known as clusters. In the one and two dimensional cases, it is easy for the human eye to quickly detect the clusters from scatter plots, assuming that the clusters exist. In higher dimensions, clustering attempts become extremely difficult without the aid of computers.

Solutions to the clustering problem usually involve the determination of a partition which satisfies some optimality criterion. The optimality criterion is a way of measuring how good a particular cluster solution is relative to other solutions. An astounding number of possible solutions exist. Reference 1 describes a Stirling number of the second kind representing the number of ways n objects may be sorted into m groups.

$$S_n^{(m)} = \frac{1}{m!} \sum_{k=0}^m (-1)^{m-k} \binom{m}{k} k^n$$

The number of groups is usually unknown so the problem is compounded, and the total number of possibilities is a sum of Stirling numbers. In the case of 25 observations, the total number of possible cluster solutions is

$$\sum_{j=1}^{25} S_{25}^{(j)}$$

which exceeds 4×10^{18} . This illustrates that the enumerative technique for finding solutions can require huge amounts of computer time, and there exists a need for a better way.

Modern techniques allow solutions to be found without evaluating the criterion for each and every solution. However the need for ranking solutions is evident, and the criterion function serves to meet this need. A wide variety of such functions exists, and the choice is usually determined by the particular characteristics of the research being conducted. A more detailed discussion of optimality criteria is presented in Section II.C.

Mathematical clustering techniques usually call for a concept of distance between objects. In order to solve the cluster problem, it is desirable to define the terms "similarity" and "difference" in a quantitative fashion. What does it mean to say two objects are different? Perhaps an investigator would assign two observations to the same group if the distance between them is sufficiently small, or to different clusters if this distance is sufficiently large. Common reference to the closeness of objects is made in

units of length, weight, or time. Numerous methods for measuring distance will be discussed in Section II.D.

In the following discussion, X_i and X_j represent two points in p -dimensional Euclidean space (E_p) corresponding to objects or observations. Any non-negative real-valued function $D(X_i, X_j)$ satisfying the following conditions qualifies as a distance function (or metric).

- a. $D(X_i, X_j) \geq 0$ for all X_i and X_j in E_p
- b. $D(X_i, X_j) = 0$ if and only if $X_i = X_j$
- c. $D(X_i, X_j) = D(X_j, X_i)$
- d. $D(X_i, X_j) \leq D(X_i, X_k) + D(X_k, X_j)$

where X_i , X_j , and X_k are any three points in E_p . Later discussions will place particular emphasis on the Mahalanobis metric.

The use of cluster analysis is applicable in nearly every field of study. The literature is both voluminous and diverse, the terminology differing from one field to another. "Numerical taxonomy" is frequently substituted for cluster analysis among biologists, botanists, and ecologists, while some social scientists may prefer "typology." Other frequently encountered terms are pattern recognition and partitioning. While discriminant analysis has been studied by statisticians for nearly 45 years, cluster analysis has only recently come to statistical notice.

Cluster analysis is an exploratory device, a tool for suggestion and discovery. A question often asked is "How do you know when you have a good set of clusters?" The answer

is that the clusters themselves are not interesting; the point of interest is in inference about the structure of the data. The clusters do not explain the structure; they are consequences of the structure. The explanatory structure is the object of the search and its description is in terms of principles and ideas, not individual data units.

It is important to realize that a given set of data may contain no "right" classification, but possibly many different, meaningful classifications. It could be the case that the data contain no clusters at all.

B. SCATTER MATRIX DECOMPOSITION

Described in this section are the multivariate terminology and notation to be used on this thesis. The literature contains as many different notational structures as authors. The emphasis is on simplicity, while also exposing the reader to some of the more common terminology.

In general, multivariate data are viewed as an n by p matrix referred to as X . It represents n observations, each of which consists of measurements on p different variables. The cross products matrix is analogous to the univariate sum of squared deviations from the mean and is represented by the p by p matrix T .

$$T = \sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - x_{..})(x_{ij} - x_{..})'$$

where

x_{ij} is the j -th observation vector in the i -th group.

$x_{..}$ is the grand mean vector of the data.

g is the number of groups.

n_i is the number of observations in the i -th group.

Prime notation indicates transpose.

All vectors are column vectors. Cross product matrices are also referred to as scatter matrices. Division of T by $n-1$ (where n represents the total number of observations) yields the total variance-covariance matrix, sometimes referred to as a dispersion matrix.

The total sum of squares (cross products) matrix may be expressed as the sum of the within-group and the between-group scatter matrices:

$$T = W + B$$

W and B are defined as follows:

$$W = \sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - x_{i.})(x_{ij} - x_{i.})'$$

$$B = \sum_{i=1}^g n_i (x_{i.} - x_{..})(x_{i.} - x_{..})'$$

where $x_{i.}$ is the mean vector of the i -th group. Each individual group has its own scatter matrix W_i , and W is the sum of these matrices:

$$W = \sum_{i=1}^g W_i$$

This discussion is intended to be completely general, with no particular group structure in mind. Later we shall explore the differences in the two group structures represented by the SOF data.

- (1) Individual SOFs are considered to be the observations and the courses are the groups.
- (2) The course mean vectors are the observations and the clusters of courses are the groups.

These two group structures are different ways of viewing the data; their relationship shall be explained in Section IV.B.

C. OPTIMALITY CRITERIA

Most of the well known clustering techniques fall into one of two main categories: (1) hierarchical and (2) partitioning. The former class is one in which every cluster obtained at any stage is a merger of clusters at previous stages. The non-hierarchical procedures however form new clusters by lumping and splitting old ones.

Partitioning methods were used in this research. The main idea is to choose some initial partition and then alter

the cluster membership in an effort to improve the partition. Different interpretations of what constitute a "better" partition and numerous ways of achieving this improvement have led to a great variety of algorithms. These methods are related to the steepest descent algorithms used for unconstrained optimization in nonlinear programming. Such algorithms begin with an initial point and then converge to a local optimum, moving one step at a time, the value of the objective function improving at each step. A well known example is the ISODATA procedure developed by Ball and Hall at Stanford Research Institute. Chapter IV discusses a partitioning method known as K-MEANS which was developed by MacQueen [2]. He uses the term "K-MEANS" to denote the process of assigning each data unit to that cluster (of k clusters) with the nearest centroid (mean vector). The cluster centroids change with each transfer of an observation.

The decomposition of the total scatter into within and between components suggests possible optimality criteria to be used in a clustering algorithm. One would like the within-groups scatter to be small relative to the between-groups scatter. Various trial clusterings could be formed using the W and B matrices as a basis for the optimality criteria which determine the best clustering. A possible choice for a criterion is to minimize trace W over all partitions into g groups. Since T is constant over all partitions, minimizing trace W is equivalent to maximizing trace B since

$$\text{trace } T = \text{trace } W + \text{trace } B$$

Although trace W is invariant under an orthogonal transformation, it is not invariant under other non-singular linear transformations.

McRae [3] points out that trace W equals the total within group sum of squares, hence the "minimum variance partition" cluster solution is found by minimizing trace W .

Considerable study has been devoted to alternative criteria such as those based on multivariate statistical analysis techniques, especially the methods of linear discriminant analysis and multivariate analysis of variance. Assuming the p variables are not linearly dependent, then as long as $p \leq n-g$, W is positive definite symmetric and so is W^{-1} . Attempts to make B and W as different as possible lead one to solving the determinantal equation:

$$|B - \lambda W| = 0$$

The solutions λ_i are the eigenvalues of the matrix $W^{-1}B$. There are t non-zero eigenvalues, where t is the minimum of p and $g-1$. This is a consequence of the fact that, if g is less than p , the g group means are contained in a $(g-1)$ -dimensional hyperplane. When $g = 2$ the analysis is equivalent to two-group discriminant analysis. Linear discriminant analysis would take the vectors originally

described in a p-dimensional coordinate system and transform the basis to a t-dimensional system. Maximizing the largest of these eigenvalues is a criterion suggested by S.N. Roy. Maximizing the trace of $W^{-1}B$, however is a criterion known as Hotelling's trace criterion. In both cases, large values for these statistics are sought in clustering algorithms since large values indicate large differences among (between) groups. Minimizing the ratio of determinants $|W| \div |T|$ is a criterion widely known as Wilks' lambda. Since T is the same for all partitions, this criterion is equivalent to minimizing $\det W$.

Both trace $W^{-1}B$ and $|T| \div |W|$ may be expressed in terms of the eigenvalues of $W^{-1}B$.

$$\frac{|T|}{|W|} = \prod_{i=1}^t (1 + \lambda_i)$$

$$\text{trace } W^{-1}B = \sum_{i=1}^t \lambda_i$$

where $t = \min(p, g-1)$. Therefore minimizing $\det W$ is equivalent to maximizing $\prod (1 + \lambda_i)$.

Friedman and Rubin [4] describe the advantages of the various criteria. Those based on multivariate statistical considerations (all but trace W) are invariant under changes in scale for the variables (non-singular linear transformation). In fact, they are the only invariants for W and B under such

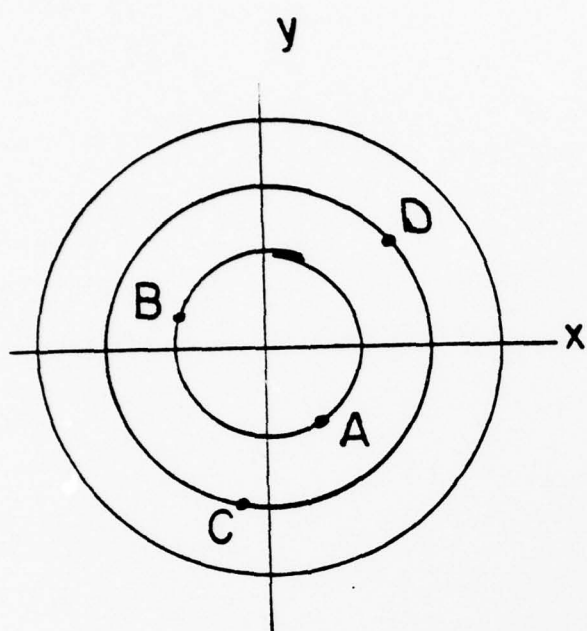
transformations. In addition, the multivariate criteria may take into account covariation among the variables.

D. DISTANCE CONSIDERATIONS

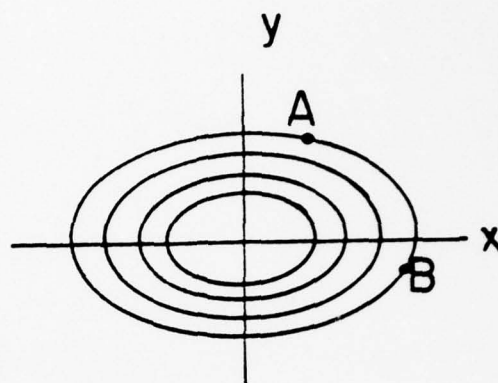
As indicated earlier there exist a number of choices for measuring distance between objects. The choice of distance function is no less important than the choice of variables to be used in the study. A serious difficulty lies in the fact that knowledge of the clusters changes the choice of distance functions. In the computation of the distance, a variable which distinguishes well between two established clusters might be weighted more heavily than others. Friedman and Rubin describe this difficulty as the "bootstrap" nature of the problem. Knowledge of the clusters would suggest an appropriate distance function which in turn would allow one to determine the original clusters. The trace W criterion implies ordinary Euclidean distance and thus hides this circularity. Use of the criteria which are invariant under non-singular linear transformations deals effectively with this circularity.

The familiar Euclidean distance is illustrated in figure 1a. When $p = 2$ the geometric interpretation of this measure amounts to determining distances by circles. Two points such as A and B on the same circle are considered equidistant from the origin, while other points such as C and D are further from the origin than A and B.

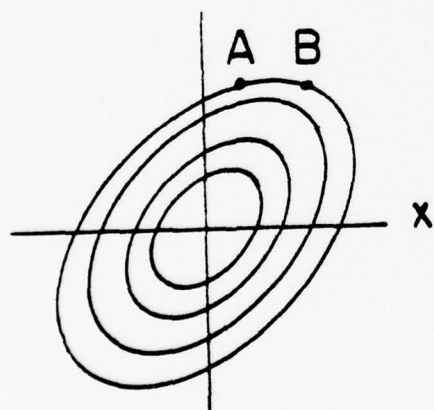
A general class of squared distance functions is provided by utilizing positive definite quadratic forms. Specifically,



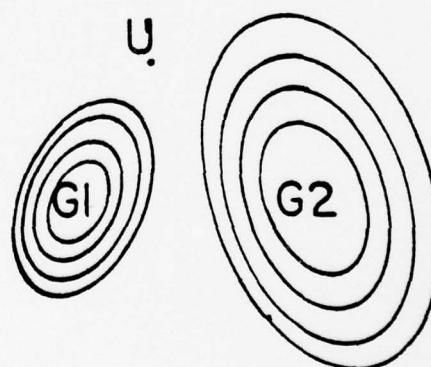
1a



1b



1c



1d

FIGURE 1

if β represents a p-dimensional observation to be assigned to one of s groups, then to measure the squared distance between β and the centroid of the i-th group one may consider the function

$$D_i = (\beta - x_{i.})^T M (\beta - x_{i.}) \quad (1)$$

where M is a positive definite matrix to ensure that $D_i \geq 0$. Different metrics are represented by different choices of the matrix M. When $M = I$ (the identity matrix) the resulting metric is the standard Euclidean distance. The variance within the data may make the unweighted Euclidean metric inappropriate. Referring to figure 1b where x has a larger variance than y, one may wish to weight a deviation in the x direction less than an equal deviation in the y direction. A method for accomplishing this is through use of an elliptical (weighted Euclidean) distance function which makes points A and B equidistant from the origin. The matrix M in this case is diagonal with diagonal elements equal to the reciprocals of the variances of the different variables. Insofar as the variance represents the true structure in the data, this distance function will adjust for differences due to the scale of measurement of each of the variables. Extending this idea further, one may consider the covariance among variables as well. Figure 1c shows how the axes may be tilted so that the major axis is oriented in a direction of reflecting the positive

correlation between x and y . Again, points on the same ellipse are considered equidistant from the origin. The matrix M in this case is the inverse of the covariance matrix.

Further examination of this concept is an important consideration in this research. If C_i represents the covariance matrix of the i -th cluster then the distance function

$$D_i = (\beta - x_{i.})^T C_i^{-1} (\beta - x_{i.})$$

uses the appropriate covariance structure when determining distance to a particular cluster centroid. Note that the number of observations in every cluster must exceed the dimensionality p in order to preserve the nonsingularity of C_i . Since C_i changes to reflect the dispersion internal to each particular cluster, the use of this metric exploits differences in the dispersion characteristics of the different groups. Figure 1d illustrates the idea. Note how a new observation (denoted by u) is closer to the centroid of group one (G_1) in terms of Euclidean distance but is more likely to be assigned to group two (G_2) when using the C_i matrix. It is instructive to point out here that if one were looking for boundaries dividing the p -dimensional space into regions, one for each of the g groups, such boundaries would be non-linear. In the performance of discriminant analysis, Eisenbeis [5] suggests appropriate quadratic classification rules.

Another choice for the M matrix in equation 1 is C^{-1} where C represents the pooled within groups covariance matrix of all the clusters.

$$C = \frac{1}{\sum_{i=1}^g (n_i - 1)} W$$

Recall from Section II.B:

$$W = \sum_{k=1}^g W_k$$

This distance is the well known Mahalanobis distance.

Note that C does not change from group to group. To ensure the non-singularity of C it must be true that $p \leq (n-g)$ where

$$n = \sum_{i=1}^g n_i$$

n represents the total number of observations over all groups.

The use of the Mahalanobis metric in the original p-dimensional space is equivalent to using Euclidean distance in the t-dimensional discriminant space with basis vectors corresponding to the eigenvectors of $W^{-1}B$. Note that the

determination of the discriminant space was based on the assumption of homogeneity of the cluster covariance structure. The Mahalanobis distance function therefore adjusts for both scale of measurement of the variables and covariation among the variables. Use of this metric is equivalent to computing distances on variables transformed to their principal components.

The natural metric to use with the trace W criterion is the Euclidean distance. However, when using criteria based on multivariate statistical considerations, Mahalanobis is the natural metric to use.

When the clusters are distributed as p-variate normal and have equal covariance matrices, then Fisher's linear discriminant function is applicable, as is the Mahalanobis distance. The accuracy of the Mahalanobis metric is sensitive to the homogeneity of the cluster dispersions and decreases as the difference between the group dispersions increases. Recall the density function for the multivariate normal distribution

$$f(x; \mu, \Sigma) = \frac{1}{(2\pi)^{p/2} (\det \Sigma)^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

where Σ is the covariance matrix and μ is the mean vector of the distribution. Note the exponent which implies utilization of Mahalanobis distance is equivalent to

measurement of the density at the point x . The empirical distributions of the clusters will therefore determine the cluster to which the observation should be assigned. The following is a proof of the invariance of Mahalanobis distance under any non-singular linear transformation. Consider the transformation

$$Y = BX$$

and let $D(Y_i, Y_j)$ represent Mahalanobis distance between Y_i and Y_j .

$$\begin{aligned} D(Y_i, Y_j) &= (Y_i - Y_j)^T C_Y^{-1} (Y_i - Y_j) \\ &= (BX_i - BX_j)^T C_Y^{-1} (BX_i - BX_j) \\ &= (X_i - X_j)^T B^T C_Y^{-1} B (X_i - X_j) \\ &= (X_i - X_j)^T B^T (BC_X B^T)^{-1} B (X_i - X_j) \\ &= (X_i - X_j)^T C_X^{-1} (X_i - X_j) \\ &= D(X_i, X_j) \end{aligned}$$

Some other common metrics are defined below.

1. L_1 norm (city block)

$$D(X_i, X_j) = \sum_{k=1}^p |x_{ki} - x_{kj}|$$

2. L_p norm (Minkowski metrics)

$$D(X_i, X_j) = \left(\sum_{k=1}^p |x_{ki} - x_{kj}|^p \right)^{1/p}$$

3. Uniform norm

$$D(X_i, X_j) = \sup_{k=1, 2, \dots, p} \{|x_{ki} - x_{kj}|\}$$

III. THE DATA SET

A. ORIGIN

The present Student Opinion Form (SOF) system was started in the summer quarter of 1975 when it replaced the Student Instruction Report (SIR) obtained from the Educational Testing Service at Princeton. The SOF form has 16 questions and space for free-form comments from the students. The information obtained from the SOF data is used for the twofold purpose mentioned in Section I.A of this paper.

A SOF form (figure 2) should be completed by each student for each course segment he takes for credit. The term "course segment" is used because the same course may be offered to more than one group of students. To differentiate between the classes, segment numbers are assigned and a separate SOF identification number exists for each segment. Different segments of the same course may or may not be taught by the same professor. About 20 percent of the forms are not returned to administration officials due to lack of interest on the part of some students and instructors. Students have been informed that the results of the SOF data are used to assist in identifying faculty members for pay raises and tenure considerations.

Difficulties with legibility of the completed forms and with the OpScan machine have persisted for several quarters. The data available for this research has been coded with

STUDENT OPINION FORM
12ND NPS 5040/215 75)

INSTRUCTOR NAME	COURSE No.	SEGMENT	NO COMMENT	STRONGLY AGREE	AGREE	NO STRONG OPINION	DISAGREE	STRONGLY DISAGREE
PLEASE USE SOFT BLACK PENCIL								
1. The course was well organized.								
2. Time in class was spent effectively								
3. The instructor seemed to know when students didn't understand the material								
4. Difficult concepts were made understandable								
5. I had confidence in the instructor's knowledge of the subject								
6. I felt free to ask questions								
7. The instructor was prepared for class								
8. The instructor's objectives for the course have been made clear								
9. The instructor made this course a worthwhile learning experience								
10. The instructor stimulated my interest in the subject area								
11. The instructor cared about student progress and did his share in helping us to learn								

PLEASE USE THE FOLLOWING SCALE FOR THE NEXT FIVE ITEMS:

0. Not Applicable / Don't know / There were none (Middle 40%)
 5. Out standing (Among the top 10%)
 4. Excellent (Among the top 30%)
 1. Poor (In the lowest 10%)
12. Overall, I would rate this instructor
 13. Overall, I would rate this course
 14. Overall, I would rate the textbook(s)
 15. Overall, I would rate the quality of the exams
 16. Overall, I would rate the laboratories

USE SPACE BELOW AND ON REVERSE SIDE FOR FREE FORM COMMENTS.
 IDENTIFY BY QUESTION NUMBER WHEN APPROPRIATE. These free form comments will be available only to the instructor.

THESE FOUR SPACES ARE AVAILABLE IF THE INSTRUCTOR WISHES TO SUPPLY ADDITIONAL QUESTIONS

PLEASE COMPLETE THE FOLLOWING ITEMS:

- ☐ THIS COURSE IS REQUIRED FOR ME
☐ THIS COURSE IS ELECTIVE FOR ME

DO NOT WRITE IN SHADED COLUMNS

QTR'S COMPT'D	QTR	HR'S THIS QTR	CURRIC	REPORT NO.
0-1	0-1	0-1	0-1	0-1
1-2	1-2	1-2	1-2	1-2
2-3	2-3	2-3	2-3	2-3
3-4	3-4	3-4	3-4	3-4
4-5	4-5	4-5	4-5	4-5
5-6	5-6	5-6	5-6	5-6
6-7	6-7	6-7	6-7	6-7
7-8	7-8	7-8	7-8	7-8
8-9	8-9	8-9	8-9	8-9
9-10	9-10	9-10	9-10	9-10

Figure 2

indications where invalid responses occur. Only the valid information was considered in this thesis. Mean scores were computed for every instructor (every course segment) from the valid responses in each of the first 13 SOF items. Only the first 13 questions were used because of the high percentage of unusable responses in items 14, 15, and 16. Each of the responses recorded is an integer from one to five, with five being the upper (more desirable) end of the scale. These data are therefore considered on an ordinal scale. Table one categorizes the blocks of data which were available for this study. Note the short 3-digit notation to be used in this paper, indicating calendar year and quarter number.

CALENDAR YEAR	NUMBER OF RESPONDENTS	3-DIGIT CODE
Summer 1977	2440	773
Fall 1977	2967	774
Winter 1978	3056	781
Spring 1978	2964	782

Table 1

The majority of the analysis was performed using only quarter 773. Unless otherwise indicated, future references to the data set shall imply quarter 773 data.

B. TRANSFORMATION

The need for a common covariance structure when using the Mahalanobis metric has been emphasized. The transformation of quarter 773 data (which attempted to accomplish homogeneity of dispersions) is explained in this section.

The SOF data are 13-dimensional, and the best transformations would involve separate examination of each of the 13 variables. Due to the overwhelming complexity of this task, only a single transformation was sought.

In the SOF data the variance is very much a function of the mean. In fact, a course with a 5.0 mean vector has no variance whatsoever. Similar effects occur on the lower end of the scale. A variance-stabilizing transformation was sought which would help to relieve the dependence of the variance on the mean. Recall the normal distribution has independent mean and variance. Other well known distributions such as the Exponential, Geometric, and Poisson all have related mean and variance. The assumption of multivariate normality underlies much of standard classical multivariate statistical methodology. The effects of departure from normality are not clearly understood. Although marginal normality does not imply joint normality, the presence of many types of non-normality is often reflected in the marginal distributions as well. The marginal distributions of the SOF data do not indicate any strong departures from normality.

Previous research by Professor R.R. Read [7] encountered the same need for a transformation of the SOF data. The

following transformation is due to Professor Read's findings:

$$\ln \left(\frac{x-1+a}{5+a-x} \right) \quad \text{where } a = .2 \quad (1)$$

The transformation was used on SOF item 12, and Bartlett's test substantiated the presence of homogeneity of variances. The groups involved here were the course segments, and the application was univariate.

Studies by Professor Glen Lindsay [8] and students in his course on Scaling Techniques produced results which suggested slight modifications to Professor Read's transformation.

$$\ln \left(\frac{x-1+a}{5+b-x} \right) \quad \text{where } a = 2.0 \quad (2) \\ b = 0.3$$

The same study could be described equally well with a constant second difference model, or what is the same thing, the function

$$x^2 + c \quad (3)$$

The three transformations were considered in the following manner. It was felt that the transformation which would produce the most nearly homogeneous covariance structure would be best. The three functions were applied to quarter

773 data, and then statistical tests for common covariance were administered. The test statistic comes from reference 5 and is explained in Appendix A. The results indicated that of the three, the first log transformation (1) generated the most nearly common covariance structure. The group structure whose covariance matrices were compared came from clusters formed by the MICKA algorithm (to be discussed in the next chapter). On the basis of the test results, the data were transformed by function (1), and all subsequent references to the data shall imply the transformed data.

Functions (1) and (2) are shown together on the graph in figure 3. The one chosen for use is the lower curve.

C. PRINCIPAL COMPONENTS

Recall the breakdown of the cross products matrix into the sum of the within and between scatter matrices. When considering the observations as individual SOFs (and the groups as courses), the cross products matrix will be called the Master scatter matrix with decomposition:

$$M = S + T$$

where S is the within course scatter and T is the between course scatter. It is reemphasized that, in this equation, the groups are the courses. The breakdown of the master scatter matrix may be examined before any clustering of course means is performed because the group structure

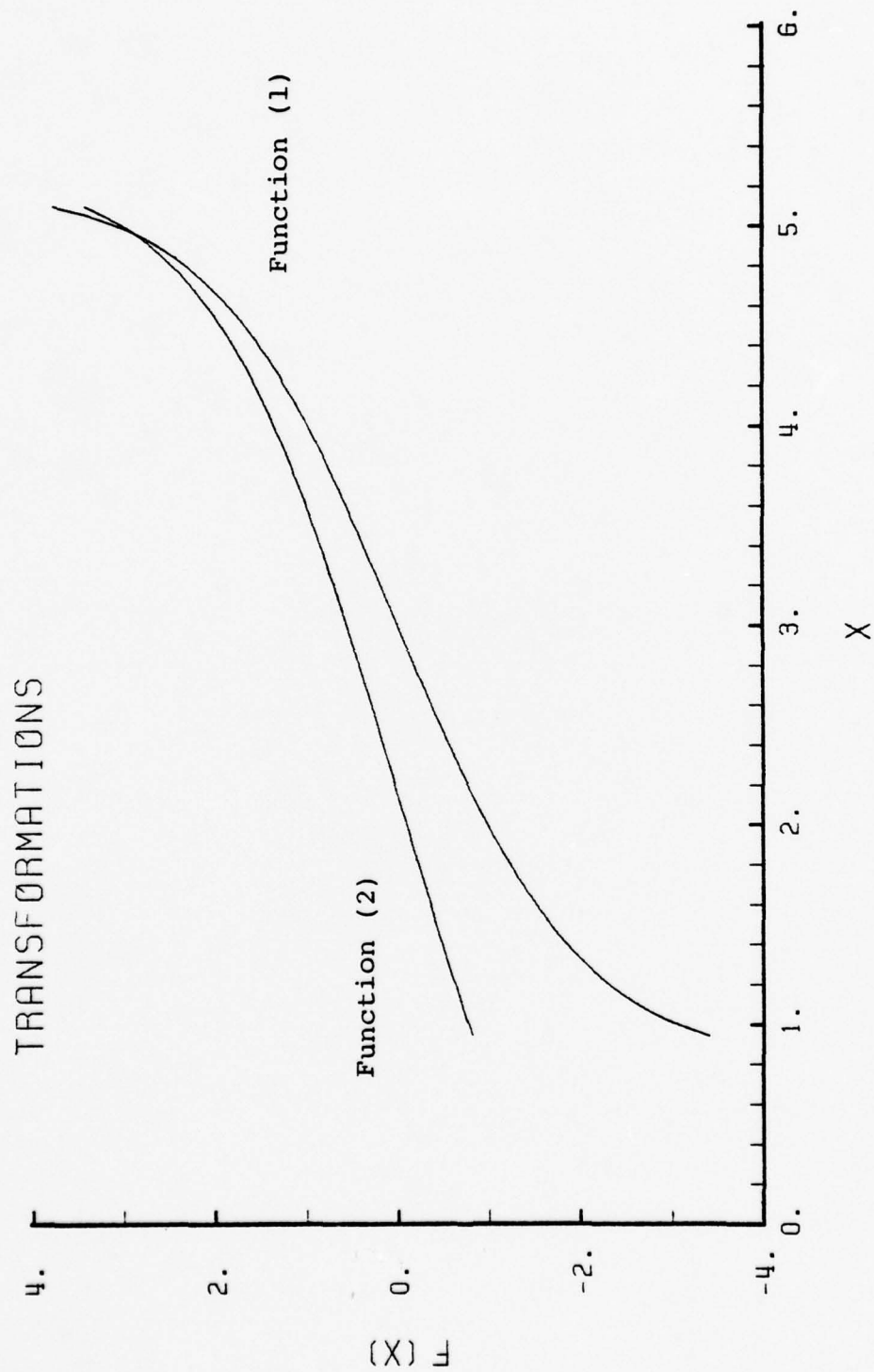


Figure 3

(courses are groups) is known. The discussion is enhanced by an algebraic description of the matrices involved. Let

$$M = \sum_{i=1}^N \sum_{j=1}^{ns_i} (X_{ij} - X_{..}) (X_{ij} - X_{..})'$$

$$S = \sum_{i=1}^N \sum_{j=1}^{ns_i} (X_{ij} - X_{i.}) (X_{ij} - X_{i.})'$$

$$T = \sum_{i=1}^N ns_i (X_{i.} - X_{..}) (X_{i.} - X_{..})'$$

where

X_{ij} is the j -th SOF response form from the i -th course.

$X_{i.}$ is the mean vector of the i -th course.

$X_{..}$ is the grand mean.

ns_i is the number of students in the i -th course.

N is the total number of courses.

Since T represents the dispersion of the course means, it is the main object of the clustering efforts. It is natural to ask also, how much information is in S . To this end a principal components analysis was performed on the covariance matrices:

$$C_T = \frac{1}{N-1} T \quad N-1 = 189$$

$$C_S = \frac{1}{\sum_{i=1} (ns_i - 1)} S \quad \begin{array}{l} \sum (ns_i - 1) = 1993 \\ \text{for quarter 773 data} \end{array}$$

Anderson [6] describes principal components as the axes of a coordinate system with special statistical properties. The principal components form a new coordinate system resulting from linear transformations of the variables which produce the special properties in terms of variances. The idea is to describe the data swarm by a new set of orthogonal coordinates so that the sample variances with respect to the new coordinates are in decreasing order. If the eigenvalues of the covariance matrix are ordered, i.e., $\lambda_1 > \lambda_2 > \dots > \lambda_p$, then the variance in the new coordinate system is greatest in the dimension associated with λ_1 , next greatest in the dimension associated with λ_2 , etc. The sum of the eigenvalues is the total variance in the original coordinate system.

The results of the principal components analysis are shown in Table 2. First, it is of interest to compute how much of the total energy in M is accounted for by T.

$$TOTAL = 18.8(1993) + 156(189) = 66952$$

PRINCIPAL COMPONENTS ANALYSIS

	C_T EIGENVALUES	C_T EIGENVECTOR FOR λ_2	C_S EIGENVALUES	C_S EIGENVECTOR FOR λ_{13}
1	0.44	0.28	0.35	-0.27
2	136.97	0.30	0.47	-0.30
3	4.64	0.28	0.49	-0.28
4	0.46	0.29	0.51	-0.29
5	3.66	0.23	0.52	-0.19
6	2.57	0.19	0.63	-0.19
7	1.14	0.25	0.63	-0.24
8	1.63	0.27	0.68	-0.29
9	1.47	0.32	0.77	-0.33
10	0.66	0.30	0.89	-0.33
11	1.23	0.26	1.00	-0.28
12	0.96	0.35	1.10	-0.30
13	<u>0.84</u>	0.26	<u>11.00</u>	-0.27
TOTAL	156		18.8	

TABLE 2

T accounts for $29484 \div 66952 = 44$ percent of the total.

This indicates that a great deal of variability must therefore be accounted for within the courses (i.e., with the students).

The principal components analysis of C_S shows the first principal component accounts for 55 percent of its total variance, but all other coordinate directions each account for 6 percent or less. Moreover, the direction of the first

component is essentially the main diagonal of 13 space, i.e., the signs are all the same and so are the magnitudes (approximately). Thus the data swarm may be thought of as an elongated ellipsoid directed along the main diagonal and having spheroidal (more or less) cross section. In particular, this suggests that the students within a course tend to score all 13 components more or less the same (all high, all moderate, or all low), but perceptions from student to student differ.

Turning to the principal components analysis of C_T , it is seen that 85 percent of the total variability is accounted for by the first principal component, and the second accounts for only three percent. Thus the data swarm of course means may be viewed as essentially one dimensional. Reference to its eigenvector reveals no single SOF item or group of SOF items is heavily weighted relative to the others and that the signs are again all the same. Thus, this component is similarly shaped along the main diagonal of 13 space, but more extremely elongated.

Some exploratory work was done on the within class variability (S) to see if the "number of quarters completed" by students has any effect on the variability represented by S. Figure four presents the results with a graph plotting within course variance versus time on board. Note the tendency for the variability to drop off in later quarters, possibly indicating more perfunctory completion of the forms.

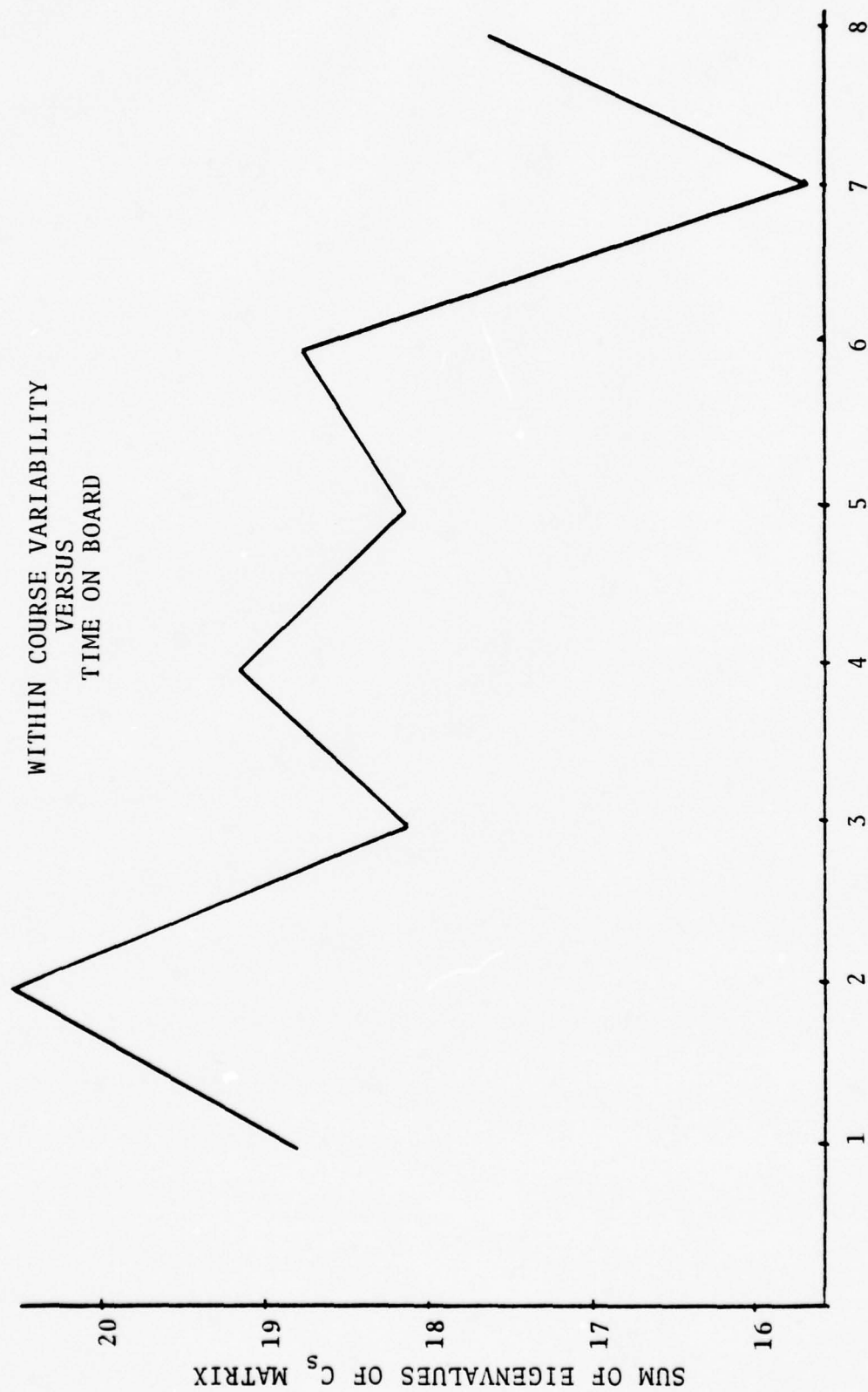


FIGURE 4

IV. THE MIKCA METHOD

A. THE ALGORITHM

The specific algorithm chosen for the cluster analysis is the MIKCA (Multivariate Iterative K-MEANS Clustering Algorithm) program written by Douglas J. McRae as a part of his doctoral dissertation at the University of North Carolina, Chapel Hill.

Reference to the flow chart in figure 5 will aid the reader in the following discussion of the algorithm. Inputs to the program are the data matrix, an estimate for g (the number of clusters), and choice of criterion and distance functions.

In the first step, preliminary calculations are made, such as the variable means and standard deviations, as well as the cross products matrix T . The next step forms the initial cluster solution. A random choice of s observations serves as the initial cluster centers. Then each of the other observations is assigned to the nearest cluster. Euclidean distance is used for this initial phase, and the cluster centroids are recomputed after each observation is assigned to a group. The observations are considered in the same order as they were input. After all of them have been assigned to clusters, the criterion value is computed. This initial cluster-finding technique is referred to as a one-pass K-MEANS procedure. It is performed three times, and

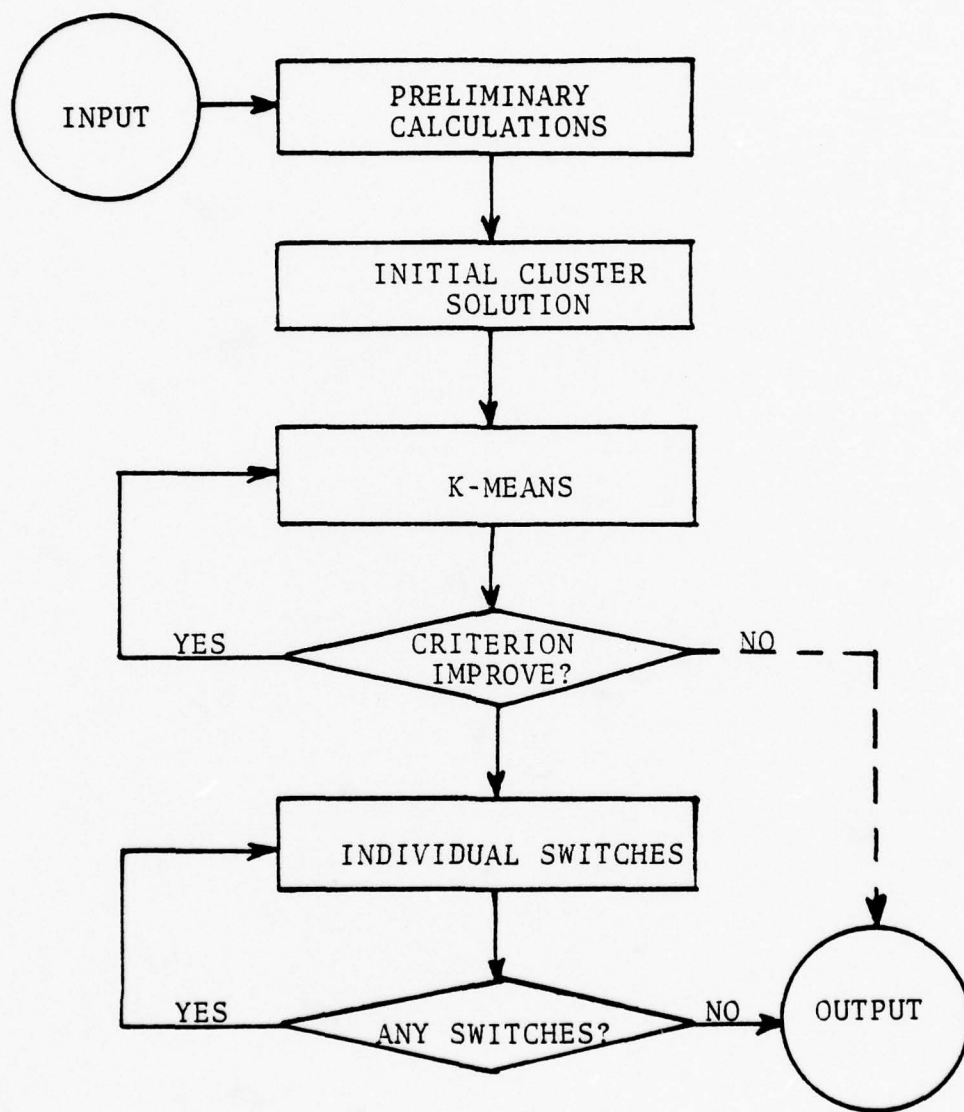


FIGURE 5 MIKCA FLOW CHART

the solution which yields the best criterion value is chosen as the initial cluster solution.

After the initial solution has been found, the program advances to the iterative K-MEANS phase where the observations are again considered in the order in which they were input to the program. It is this phase where the user's choice of distance function is used. The distance from each observation to each cluster centroid is again computed, this time with the user's distance function, the assignment to the closest centroid being made and the centroid updated to reflect its new membership. After considering all n observations in this manner, the new criterion value is checked for possible improvement during the K-MEANS iteration. As long as the criterion value improves, the K-MEANS procedure is repeated; if the criterion fails to improve then the MIKCA algorithm goes to the next step, the individual switches section.

Note the importance of the order of consideration of the observations. The order is important because the cluster means are recomputed after each observation is reassigned.

In the individual switches phase, consideration is given to moving each observation to every other cluster, the move being made if and only if an improvement in the value of the criterion results. An elaborate labelling procedure provides a unique order in which to consider each observation. This procedure continues until a complete pass through the data is made with no changes in cluster membership.

The MIKCA algorithm provides the following options for distance and criterion functions.

CRITERION

1. Minimum trace W
2. Minimum $\det W$
3. Maximum largest order of $|B - \lambda W| = 0$
4. Maximum sum of roots of $|B - \lambda W| = 0$

DISTANCE

1. Euclidean
2. Weighted Euclidean
3. Mahalanobis

Using R.A. Fisher's iris data, McRae tested his algorithm and produced extremely good results. Using the $\det W$ criterion and Mahalanobis distance, MIKCA produced a solution identical to the classification given by multiple discriminant analysis. This is a notable achievement since the cluster procedure, which does not know the true composition before the analysis, makes the same final classification of observations as does the discriminant procedure, which bases its analysis on the group composition information.

The MIKCA provides as output the value of the criterion function, the cluster membership, and the cluster mean vectors. Also provided are two matrices, T and W . The program was written in FORTRAN IV for the IBM 360 series of computers.

B. MODIFIED MIKCA

Initially, the MIKCA program was used with the p-component mean responses for each course as the input data matrix. Since the number of students utilized in producing these means is quite variable, these input vectors are not equally well determined and, as has been mentioned earlier, may effect the covariance structure between the objects. It is desirable to have the option of weighting these course means in order to effect a better balance in terms of their accuracy and to reduce any consequential distortion in the covariances. It is convenient to refer to this modification as the "1 man 1 vote" option, and to the original technique as the "1 course 1 vote" option. The following algebraic definitions will aid in illustrating the weighting effect.

Recall the breakdown of the master scatter matrix into the sum of within and between matrices.

$$M = S + T$$

When the mean scores are computed for each course and used as inputs to MIKCA, then a different dispersion structure takes form. The groups are no longer the known courses, but are now the object of the problem. The groups are unknown clusters of courses (or professors). Let T^* denote the total scatter contained in the data when each observation represents a course mean vector. T^* may also be expressed as the sum of within and between scatter matrices.

$$T^* = W^* + B^*$$

These matrices are defined as follows:

$$T^* = \sum_{s=1}^g \sum_{k=1}^{nc_s} (\bar{x}_{sk} - \bar{x}_{..}) (\bar{x}_{sk} - \bar{x}_{..})'$$

$$W^* = \sum_{s=1}^g \sum_{k=1}^{nc_s} (\bar{x}_{sk} - \bar{x}_{s.}) (\bar{x}_{sk} - \bar{x}_{s.})'$$

$$B^* = \sum_{s=1}^g nc_s (\bar{x}_{s.} - \bar{x}_{..}) (\bar{x}_{s.} - \bar{x}_{..})'$$

where

nc_s is the number of observations (courses) in the s -th cluster.

g is the number of clusters.

\bar{x}_{sk} is the k -th observation (course mean vector) in the s -th cluster

$\bar{x}_{s.}$ is the mean vector of the s -th cluster

$$\left(\frac{\sum_k \bar{x}_{sk}}{nc_s} \right)$$

$\bar{x}_{..}$ is the grand mean

$$\left(\frac{\sum_s \sum_k \bar{x}_{sk}}{\sum_s nc_s} \right)$$

Note that the grand mean mentioned here is not the same as the grand mean used in the decomposition of the master matrix M. The difference between T and T* is that T is weighted by the number of students in each course, ns_i . This weighting factor was lost when the individual observations were viewed as the class mean vectors. A close algebraic examination of T will illustrate its weighted property. Originally, we had $M = S + T$ where:

$$T = \sum_{i=1}^N ns_i (x_{i.} - x_{..}) (x_{i.} - x_{..})'$$

It is now helpful to show the decomposition of T.

$$T = W + B$$

Let $x_{i.}$ become \bar{x}_{sk} (k-th course mean in s-th cluster) and ns_i become ns_{sk} (number of students in k-th course of s-th cluster). Therefore the same T can be reexpressed as

$$T = \sum_{s=1}^g \sum_{k=1}^{nc_s} ns_{sk} (\bar{x}_{sk} - x_{..}) (\bar{x}_{sk} - x_{..})'$$

Letting

$$W_s = \sum_{k=1}^{nc_s} ns_{sk} (\bar{x}_{sk} - \bar{x}_s) (\bar{x}_{sk} - \bar{x}_s)'$$

where

$$\bar{x}_s = \frac{\sum_{k=1}^{nc_s} ns_{sk} \bar{x}_{sk}}{\sum_{k=1}^{nc_s} ns_{sk}} \quad (\text{weighted mean vector of } s\text{-th cluster})$$

then

$$W = \sum_{s=1}^g W_s$$

and

$$T = W + B \quad (B \text{ is obtained by subtraction})$$

The understanding of this distinction is important because it describes the abbreviated (unweighted) dispersion upon which MIKCA bases its cluster solution.

A number of changes were made to the MIKCA computer code to allow for a system of weights, ns_i , for the course means. The modified code extends the capability of MIKCA by making this option available. It amounts to using T rather than T^* as the basic dispersion structure. This seems more natural because the matrix T appears in the earlier decomposition.

$$M = S + T = S + W + B$$

Some of the changes are summarized here:

1. Allow for class size as input.
2. Alter the computation of T to allow for the weighting factor.
3. Alter the computation of cluster centroids to allow for weighting.
4. Alter calculations of the B matrix for the same reason. (W is found by subtracting B from T .)

The computer code for the modified MIKCA is included in Appendix F.

Cluster solutions using both weighted (T) and unweighted (T^*) dispersion structures were found and compared (see table 3 in next section). The comparison indicates some differences in cluster solutions, however the importance of these differences is left to the reader.

C. RESULTS

Several cluster solutions were formed using the MIKCA algorithm. It seemed wise to include the number of students in a course as the 14-th variable. The natural logarithm of the class size was the transformation applied to this variable. Since class sizes ranged from two to 40, this transformation brought the values into a similar range as the other 13 variables and also reduced skewness. For quarter 773 the mean class size was 12.7 students with a standard deviation of 7.9. For the transformed variable these

statistics are 2.3 and 0.7. Cluster solutions were found with and without inclusion of this 14-th variable. The results are shown in table three.

Another option available to the MIKCA user is the standardization of variables prior to entering the clustering process. McRae points out how this option becomes very useful when the variables are on vastly different scales of measurement. Except for the 14-th variable the present scales are psychological in nature and seem to be much the same. Some exploratory work was performed with the standardization option (see table three) but it was not considered significant because of the similarity in the scales of measurement.

Table three shows the comparisons of cluster results obtained under various conditions. The comparison coefficient provides a measure of agreement between solutions and is computed by a method introduced in Chapter VII. Table three shows generally higher values for $g = 3$, indicating that there exists robustness of solutions for the smaller values of g .

The results of these cluster solutions may also be seen in graphical form by referring to Appendix B. These graphs, called profile charts, depict the mean vectors for each of the clusters formed by the MIKCA algorithm. The mean vectors have been standardized so that one can see the number of standard deviations from the grand mean. These profiles are .

COMPARISON COEFFICIENTS FOR CLUSTER SOLUTIONS
OBTAINED UNDER VARIOUS CONDITIONS

13N	13S	14N	14S	13N	13S	14N	14S	13N	13S	14N	14S
1.0	1.0	.75	.83	1.0	1.0	.63	.61	1.0	.57	.72	.82
	1.0	.75	.83		1.0	.63	.61		1.0	.44	.48
		1.0	.86			1.0	.58			1.0	.75
			1.0				1.0				1.0

LEGEND

13N 13 VARIABLES, NOT STANDARDIZED
13S 13 VARIABLES, STANDARDIZED
14N 14 VARIABLES, NOT STANDARDIZED
14S 14 VARIABLES, STANDARDIZED

COMPARISON OF
ORIGINAL MIKCA
AND
MODIFIED MIKCA
SOLUTIONS

COMPARISON COEFFICIENT = .57

TABLE 3

also helpful in identifying the variables which are significant in the cluster membership. For example, an important variable would be one that produces a break in the pattern.

In the 13 variable case, the profiles produced results which indicated the lack of clearly dominant variables in cluster identification. With introduction of the 14-th variable, some very revealing results become immediately apparent. While the cluster membership changed little in going from 13 to 14 variables, the cluster with the highest mean vector became clearly associated with the smallest class sizes. Similarly the cluster with the lowest mean vector is characterized by a very large class size. This finding is one of the most significant results.

One of the most critical decisions facing the analyst is the number of clusters to form. Some algorithms based on the K-MEANS idea allow g to change during the clustering process, however the MIKCA method requires g to be input by the user and it does not change in the course of the program execution. Typically the investigator does not know the number of clusters in the data, and he must make some educated guess. As pointed out earlier, it is possible for several different, but meaningful, cluster solutions to exist in one body of data.

The method used to determine g was to obtain solutions based upon several values for g and then plot the criterion values for each of these solutions. An appropriate choice for g would be a number beyond which the marginal improvement

of the criterion becomes insignificant. Figures 6 and 7 are the results of such tests suggesting that six clusters represent the major portion of the separating power of the algorithm.

Profile charts of the cluster solutions with $g = 6$ were uninteresting. The middle clusters were all bunched together suggesting that clusters were forced on that part of that data where perhaps they did not actually exist (i.e., sparse data near the boundaries). Comparison results (table 3) indicate a much more stable solution when g is reduced below six.

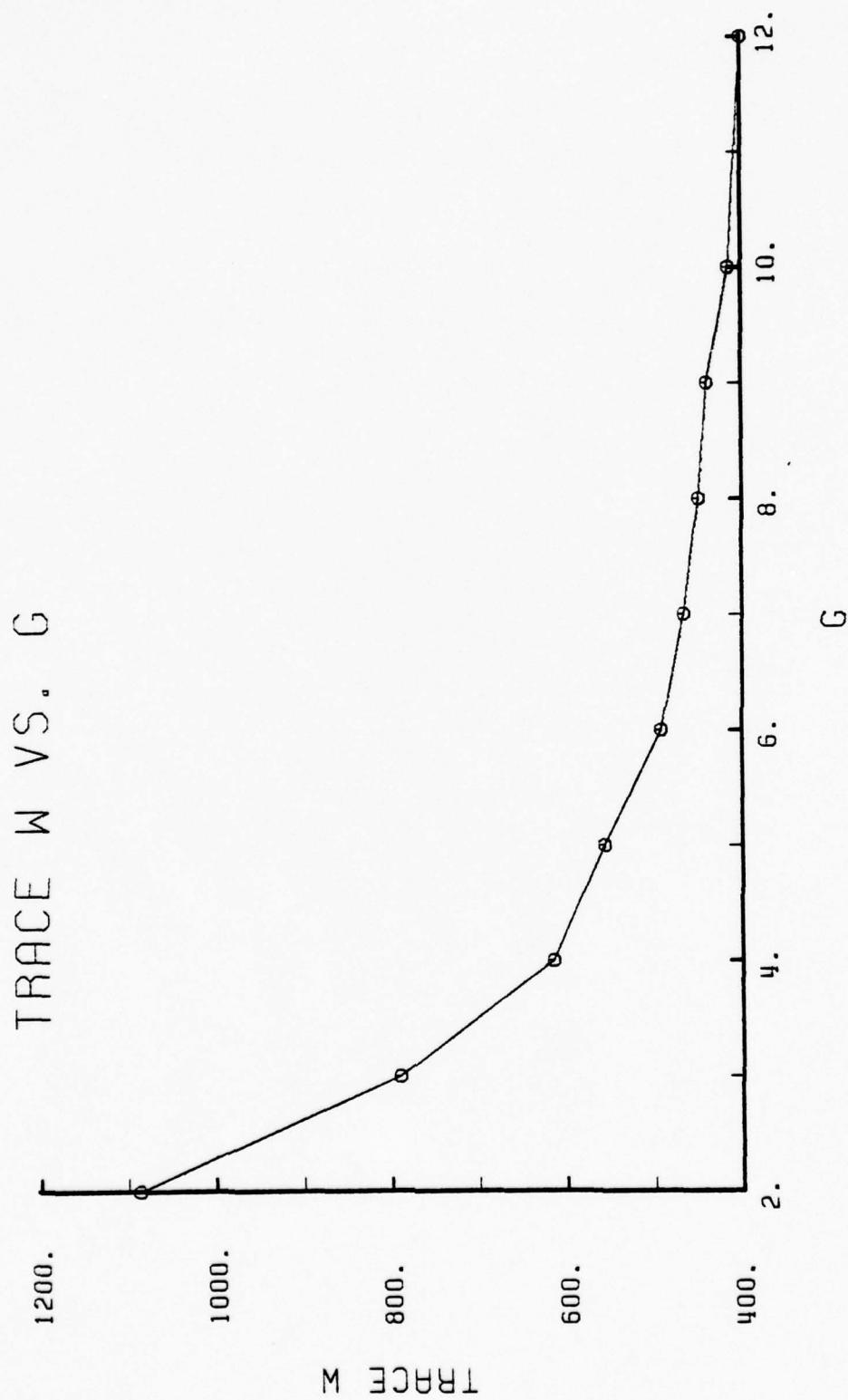


Figure 6

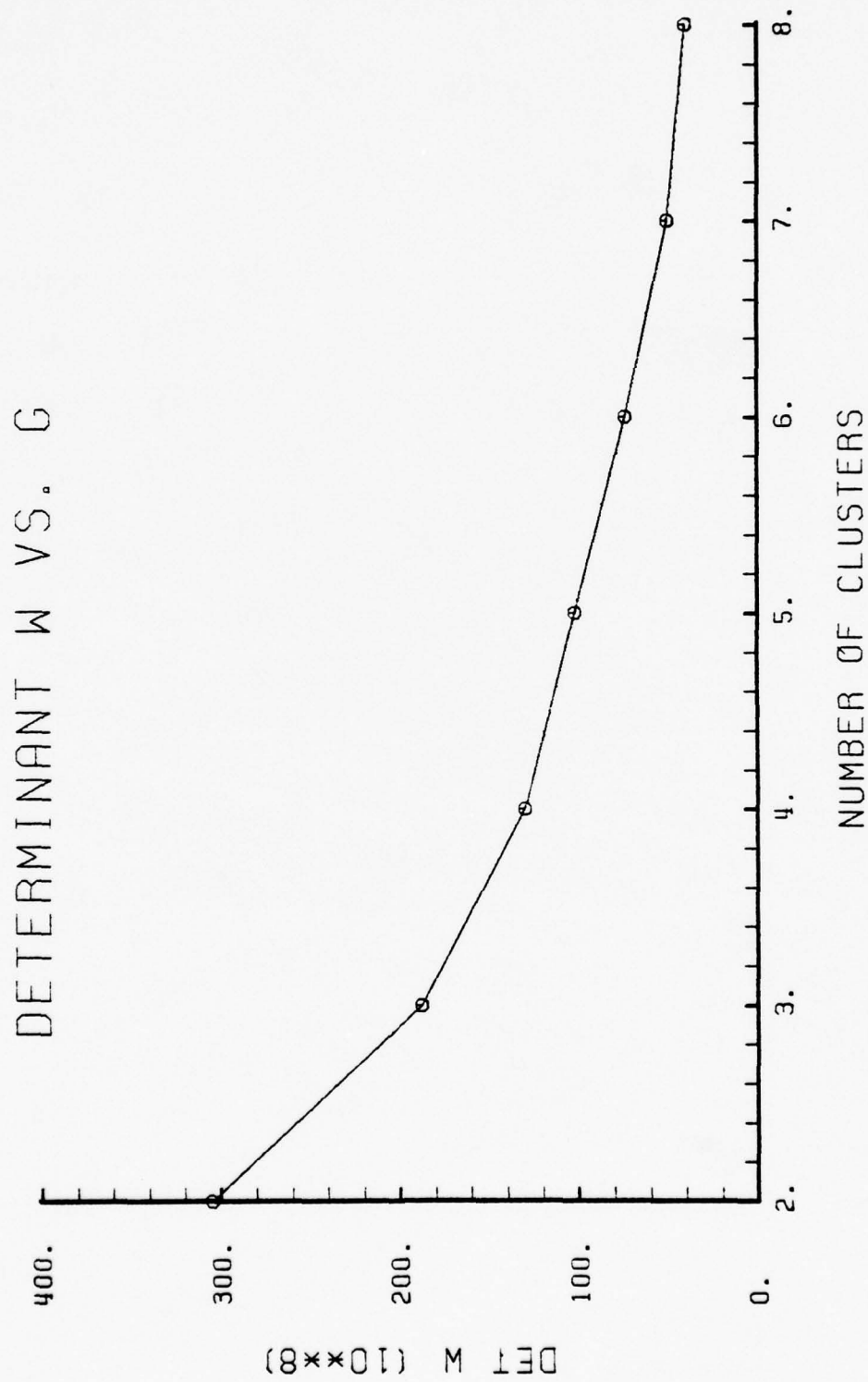


Figure 7

V. DISCRIMINANT ANALYSIS

A. THEORY

As mentioned earlier, discriminant analysis allows an analyst to classify new observations based on observations which are samples from known groups. Only Fisher's linear discriminant function will be used in this study. It also provides information about the relative importance of the various variables in assigning an observation to a group. The linear discriminant function is based on the assumptions of multivariate normality and homogeneity of dispersions. The ability to identify the dominant variables and the dimension reduction offered by the discriminant space were both extremely useful aids for analyzing the SOF data.

These "more important" variables will be earmarked for later use in the construction of Chernoff's FACES. Also of interest is the plot of data points in discriminant space. The interaction of the coefficients in the discriminant functions will be seen as well as the characterization of the dimensions.

In order to describe our usage of discriminant analysis, let us first suppose there are only two clusters in 13-dimensional space. It is desired to project these two clusters orthogonally onto a line so that the variation between the two groups is as large as possible relative to the variation within the two groups. Finding the direction of projection to accomplish this is part of the purpose of

discriminant analysis. The solution provides a way of discriminating between the two clusters by a suitable linear combination of the 13 variables. The same theory is generalized to g groups, where Wilks [9] has shown that a projection to the smaller of $g-1$ or p dimensions is possible without loss of information. Recall the earlier discussion that indicated this smaller number as t , the number of non-zero eigenvalues of $W^{-1}B$. The eigenvalues are the variances in the direction of their associated eigenvectors. One can easily determine the proportion of variance attributable to each of the dimensions of discriminant space and also the SOF items which load most heavily in each dimension.

One gains insight into the variables from examination of the coefficients in the discriminant functions. There is one function for each dimension, the standardized coefficients of which are used in this analysis.

B. RESULTS

Up to this point, most of the analysis has been performed on the 190 courses in quarter 773. A smaller, more manageable data base was needed to continue. Also, it seemed wise to prepare to study individual departments. The Electrical Engineering Department was chosen for further analysis since it is a large department and hence not too small for this purpose. Over the four quarter period, there were 116 course segments with valid SOF responses. These 116 courses from the EE department were the data used in the discriminant analysis.

When dealing with four clusters, the dimensionality of the discriminant space is three, and depending on the size of the eigenvalues, perhaps fewer dimensions will provide sufficient discrimination. Table four gives the results of performing a discriminant analysis. Figure eight is a graph of the two-dimensional discriminant space (the third dimension is neglected).

The eigenvalues indicate 94 percent of the total variance is represented by the first two discriminant functions. Figure eight corroborates this fact by depicting easily seen separation in two dimensions. Imagine projecting the points onto the horizontal axis. Discrimination in the first dimension would account for 73.6 percent of the variation. Groups one and four would easily be separated, but two and three would overlap.

Examination of the coefficients will enable one to label the dimensions by identifying the dominant characteristics which they measure. The first dimension is along the horizontal axis and is associated with the first discriminant function. The magnitude of the coefficients indicates their relative impact on the dimension. The signs aid in understanding which variables reinforce one another (matching signs) and which tend to cancel (opposite signs). In the first function of table four, SOF item 12 is the most prominent. This question (see figure 2) asks the student to score the overall rating of the instructor. It is not surprising

RESULTS OF DISCRIMINANT ANALYSIS ON THE 116 COURSES
IN EE DEPARTMENT OVER A FOUR QUARTER PERIOD

DISCRIMINANT FUNCTION	EIGENVALUE	RELATIVE PERCENTAGE
1	5.79	73.6
2	1.64	20.8
3	0.44	5.6

STANDARDIZED DISCRIMINANT
FUNCTION COEFFICIENTS

	<u>Function 1</u>	<u>Function 2</u>	<u>Function 3</u>
1	-0.11	0.23	-0.22
2	0.13	0.14	-0.09
3	-0.05	1.46	0.01
4	-0.47	-0.36	-0.80
5	-0.31	0.01	-0.82
6	0.08	0.36	-0.74
7	0.15	-0.36	0.12
8	0.23	1.15	-0.36
9	0.36	-0.82	-0.47
10	0.05	0.08	-0.77
11	-0.20	-1.16	1.18
12	-0.72	-1.08	1.80
13	-0.18	0.91	0.92

TABLE 4

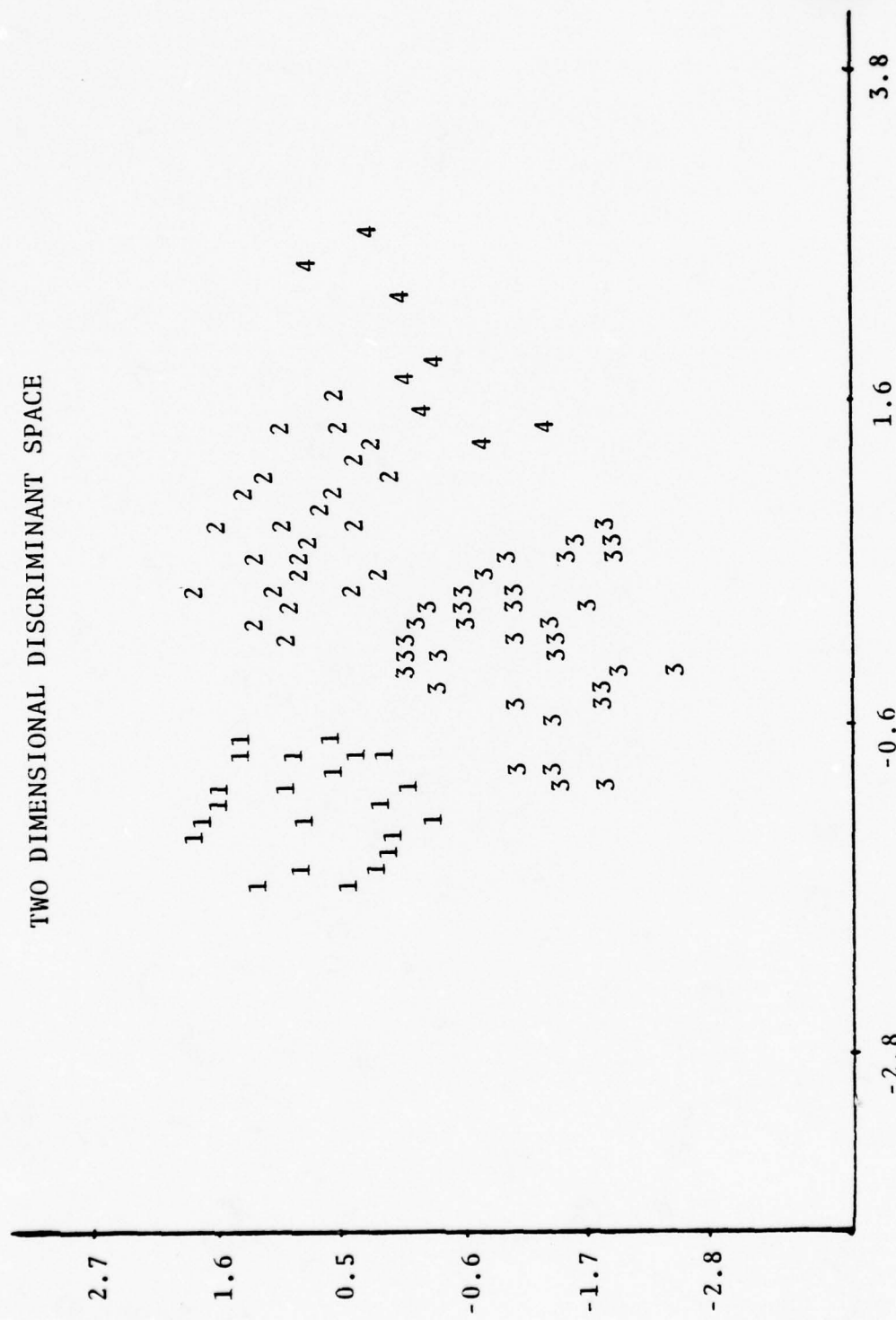


FIGURE 8

that this question is very important in the discrimination process. High marks on items four and five tend to reinforce a high mark on question 12. Those questions are:

- (4) Difficult concepts were made understandable.
- (5) I had confidence in the instructor's knowledge of the subject.

Interestingly, a high mark on item nine (instructor made the course a worthwhile learning experience) tends to diminish the effect of a high score on item 12. The first dimension is dominated by question 12 and was labeled the "popularity" dimension.

The second dimension is depicted by the vertical axis on the graph in figure eight, and measurements along this dimension are controlled by the second discriminant function. The separating power in this direction is less than one third that of the first. Note however that the vertical scale is compressed 25 percent more than the horizontal scale (1.5 inches vertical = 2.0 inches horizontal). Items three and eight has strong positive coefficients whereas questions 11 and 12 are pulling heavily in the negative direction. However the strength of the information is not great, and deeper interpretation hardly seems worth the effort.

Only 5.6 percent of the total variance appears in the third function, and it is therefore considered insignificant. One might note that item 12 also dominates the third dimension.

The main purpose here has been to identify variables for use in constructing Chernoff FACES. The discriminant analysis has served that purpose well, and it has also described the character of the dimensions.

VI. CHERNOFF FACES

A. BACKGROUND

Chernoff's FACES was the second cluster method to be applied to the SOF data. The method was used with the same purpose in mind, and it was hoped that earlier cluster solutions could be reproduced by this method. Additionally, there was the possibility of gaining new information about the structure within the data. Professor Herman Chernoff developed this graphical method for representing multivariate data. The now familiar data point in p-space is represented by a computer-drawn cartoon of a face whose characteristics (features) are determined by the position of the point. Features such as nose length and mouth curvature correspond to components of the data point. In the case of the SOF data, each component of the 13-dimensional vector can be made to control one of 20 features, and seven constants can be selected for the remaining features. The technique lends itself to clustering since the investigator can group together those faces which resemble each other.

Chernoff [10] points out that people spend a great deal of their life studying and reacting to faces. The human mind subconsciously acts as a high speed computer sometimes detecting barely measurable differences and ignoring unimportant differences, even if they are large. Chernoff claims that unlike a machine, the mind has the capability

to disregard non-informative data and search for useful information. He states that certain major characteristics of the faces are instantly observed and easily remembered; finer details and correlations become apparent after studying the faces. Clustering by sorting the faces is certainly easier than staring at a large matrix of data. The method has pitfalls and limitations and some of them will be dealt with in this thesis.

After the publication of Chernoff's method [11], quite a number of people began experimenting with the technique. Lake [12] mentions a few more successful applications of Chernoff's method, including:

1. L.A. Bruckner of Los Alamos Scientific Lab of the University of California while studying the performance of offshore oil companies.
2. Johns Hopkins University
 - a. Developing methods of psychiatric screening.
 - b. Monitoring patients in intensive care units.
 - c. Monitoring the stock market.
3. Dr. David L. Huff of the University of Texas in developing urban regional indicators that measure the quality of life.
4. Professor P.C.C. Wang and Gerald Lake at the Naval Postgraduate School in analyzing Soviet naval penetrations into the Indian Ocean and the African littoral; and Soviet foreign policy in sub-Saharan African states.

5. Professor Chernoff in geological and fossile-related experiments.

The field of computer graphics has experienced tremendous growth in recent years due mainly to the state of the art in computers and computer display equipment (including both video and plotting types). The adage that "a picture is worth a thousand words" has proven to be quite true. Recent developments include on-line programs that perform statistical analysis with polygon, bar graphs, arrows, and scatter diagrams. Three-dimensional data displays have facilitated the work of engineers and statisticians alike.

An interesting application of the FACES program is Bruckner's study of offshore drilling by oil companies. Figure 9 displays some of his results. Two of the features, nose width and eye separation, are controlled by the variables "expected years to production" and "number of leases won", respectively. Other features are controlled by a variety of variables representing the company's financial health and growth potential.

Reference to figure 10 will help describe how the faces are constructed. Table 5 gives the range of the variables which control the features and distance parameters of the face.

The data are first converted to the X parameters as follows. The variable Z is used to control the parameter X_i which is allowed to range from a_i to b_i .

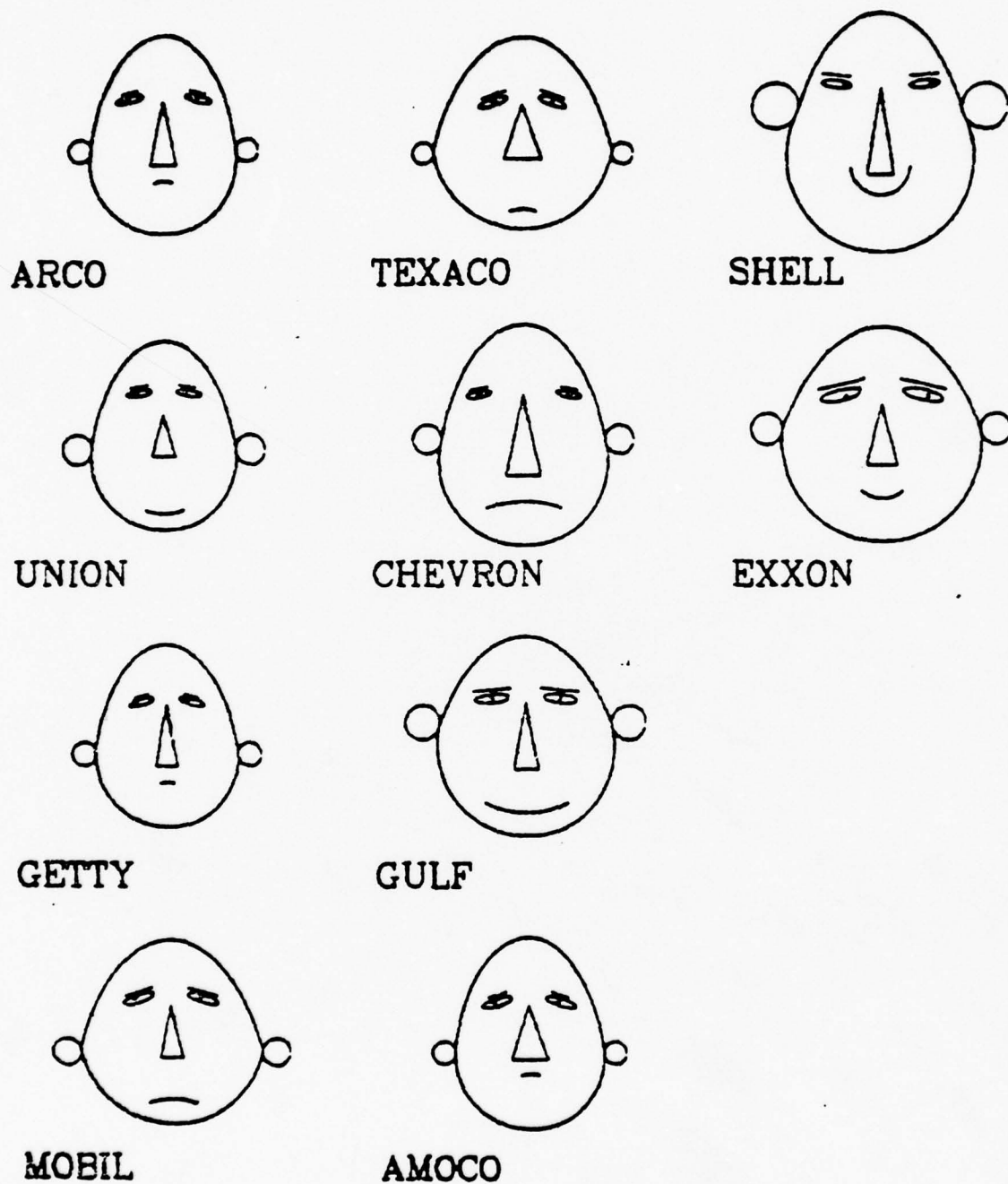


Figure 9

Bruckner's Offshore Hydrocarbon Producers

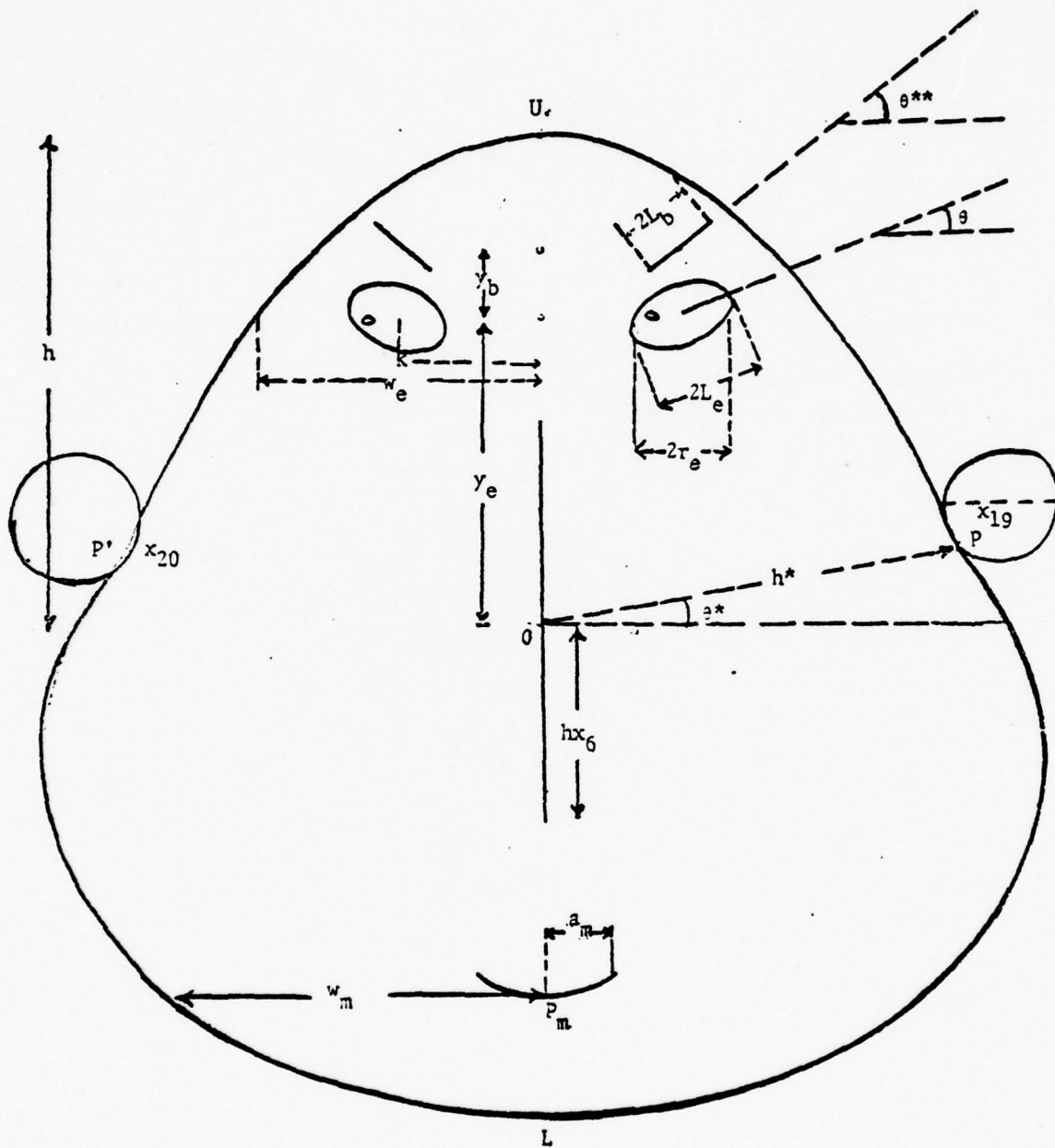


Figure 10

Chernoff Face with Ears

FEATURE RANGES AND DESCRIPTION

This table is taken from reference 10, and the descriptions are not complete. For a more detailed, mathematical explanation, see Appendix C.

<u>Range</u>			
(0,1)	x_1	controls h^*	distance from 0 to P
(0,1)	x_2	controls θ^*	angle between OP and horizontal
(0,1)	x_3	controls h	half-height of face
(0.5,2)	x_4	is	eccentricity of upper ellipse of face (width/height)
(0.5,2)	x_5	is	eccentricity of lower ellipse of face (width/height)
(0,1)	x_6	controls	length of nose
(0,1)	x_7	controls P_m	position of center of mouth
(-5,5)	x_8	controls	curvature of mouth (radius = h/x_8)
(0,1)	x_9	controls a_m	length of mouth
(0,1)	x_{10}	controls y_e	height of centers of eyes
(0,1)	x_{11}	controls x_e	separation of centers of eyes
(0,1)	x_{12}	controls θ	slant of eyes
(0.4,0.8)	x_{13}	is	eccentricity of eyes (height/width)
(0,1)	x_{14}	controls L_e	half-length of eye (L_e also depends in part on x_{10} and x_{11})
(0,1)	x_{15}	controls	position of pupils
(0,1)	x_{16}	controls y_b	height of eyebrow center relative to eye
(0,1)	x_{17}	controls $\theta^{**}-\theta$	angle of brow relative to eye
(0,1)	x_{18}	controls	length of brow
(0,1)	x_{19}	is	ear diameter
(0,1)	x_{20}	is	nose width

TABLE 5

$$X_i = a_i + (b_i - a_i) \left(\frac{Z - m}{M - m} \right)$$

where m and M are the observed minimum and maximum of Z .

Chernoff's technical report [10] presents a very detailed description of the geometric relationship of the features in the face construction. A few general remarks concerning the geometric attributes are included here. The boundary of the face is formed by joining portions of two ellipses, an upper and a lower. The angle θ determines where the ellipses meet and consequently, the height of the ears. The nose is a triangle centered at the origin. Both its height and width are variable. The curvature of the mouth is a portion of a circle, the radius and center of which are also variable. The eyes are formed by ellipses whose angle, half-length, and eccentricity are all controlled by variables.

B. FEATURE-VARIABLE RELATIONSHIP

A frequent question is whether some features are more informative than others. Some observers feel that the eyes convey the most information; others regard the mouth or the shape of the face as the most relevant feature. The results of the discriminant analysis identified the most dominant variables in the discriminant space. Now these variables must be assigned to facial features.

Chernoff [13] himself conducted an experiment to evaluate the effect on classification error of random permutations in the assignment of variables to features. He found that

random permutations would change the faces so that a classifier might increase or decrease his number of errors by a factor of about 25 percent. Unfortunately, his experiment did not evaluate the efficiency of specific features. His studies also make no effort to determine whether ability to discriminate depends on the dimensionality of the data.

Considering Chernoff's findings, it would seem that the assignment of variables to features is of minor importance. The use of discriminant analysis provides a way of detecting which variables are important, and it seems appropriate to take advantage of this valuable information when constructing the faces. Moreover, there is choice in the features that are selected for use. The author's choice of the six best features are starred in table 6. The table gives the complete list of feature-variable combinations. The results of the discriminant analysis were relied upon heavily in forming the variable assignments.

Reference to figure eight (discriminant space) and table 6 will aid in the following discussion. In the first dimension the important SOF items are 12 and 4 which control the mouth curvature and ear height, respectively. High scores on these two items separate the observation well to the negative end of the scale and cause the face to have a big smile and high ears. Items 12 and 4 have the same sign (negative) but item 9 is associated with a large positive coefficient and controls the lower eccentricity of

FEATURE-VARIABLE COMBINATIONS

FEATURE	THREE DIFFERENT TRIALS CONTROLLED BY		
	13 VAR	6 VAR	8 VAR
1 FACE WIDTH	0.5	0.5	0.5
2 ANGLE θ	4	0.65	4
3 FACE HEIGHT	0.7	0.7	0.7
4 UPPER ECCENTRICITY	8	0.95	0.95
*5 LOWER ECCENTRICITY	9	4	0.6
6 NOSE LENGTH	10	0.45	9
7 MOUTH CENTER	0.5	0.3	0.5
*8 MOUTH CURVATURE	12	12	12
9 MOUTH LENGTH	13	0.7	0.8
10 EYE HEIGHT	0.23	0.23	0.23
11 EYE SEPARATION	1	0.5	0.5
*12 EYE SLANT	11	3	3
13 EYE ECCENTRICITY	3	0.6	0.6
14 EYE HALF LENGTH	6	0.5	5
*15 PUPIL POSITION	2	9	13
16 EYEBROW HEIGHT	0.3	0.3	0.3
*17 EYEBROW ANGLE	5	8	8
18 EYEBROW LENGTH	0.4	0.4	0.4
19 EAR DIAMETER	0.3	0.3	0.3
*20 NOSE WIDTH	7	11	11

Integer numbers are the SOF item #.
 Decimal values are the fixed features

TABLE 6

the face. A low mark on this item would complement high marks on items 12 and 4 and would be reflected in the lower face having small eccentricity (more narrow).

Turning to the vertical axis (second dimension) of figure eight which has 21 percent of the total variance, the dominant variables are 3, 8, and 11, where 11 is negative; 3 and 8 are positive. High scores on items 3 and 8 separate the observation upward on the vertical axis and are reflected as highly eccentric eyes and upper ellipse.

Droopy eyes, reflecting a small value for SOF item 11, tend to complement and reinforce the higher values for items three and eight. It seems like a good idea to use the results of the discriminant analysis in this way, but it is impossible for the viewer to know which variables act together and which interfere unless he is told beforehand.

A good deal of exploratory work was carried out to determine useful ranges for the features. The more the features are allowed to vary, the wider the variety of faces produced. With large ranges, however, faces formed from extreme data can become very distorted. On the other extreme, too little variability in the ranges suppresses valuable information and hinders the clustering process.

It appears that the best ranges depend on the structure and amount of variability in the data. Every data set has its own characteristics, and it is best to tailor each to its own best set of ranges. A great portion of the SOF data is found "close" to the grand mean, but with a few

significant outliers. In order to provide discriminating ability among the largest mass of the data, the appearance of the outliers was further accentuated. The ranges were set at values which would allow close-in discrimination, but simultaneously attempted to minimize the departure of the outliers.

C. CLUSTERING THE FACES

The next step in this research is to cluster the faces. This task was performed by six students in the Operations Research curriculum. The faces are shown in figure 11; the 33 course segments from the Electrical Engineering department in quarter 781. The judges (students performing the clustering) were given no information concerning the feature-variable combination. They were simply instructed to group the faces in the manner which best suited them. Fifteen minutes were allowed for the task. The purpose was to quickly, but carefully, cluster the faces. The judges were reminded that each face is different and to search for the most natural looking clusters. It was felt that too much time spent on this task could defeat the purpose of the faces as a first pass look at the data. In every case, the judges acted independently of one another. No clues were provided which might have indicated which features were more important.

Figure 11 shows the faces in the clusters which were formed by the MIKCA algorithm. This cluster structure was

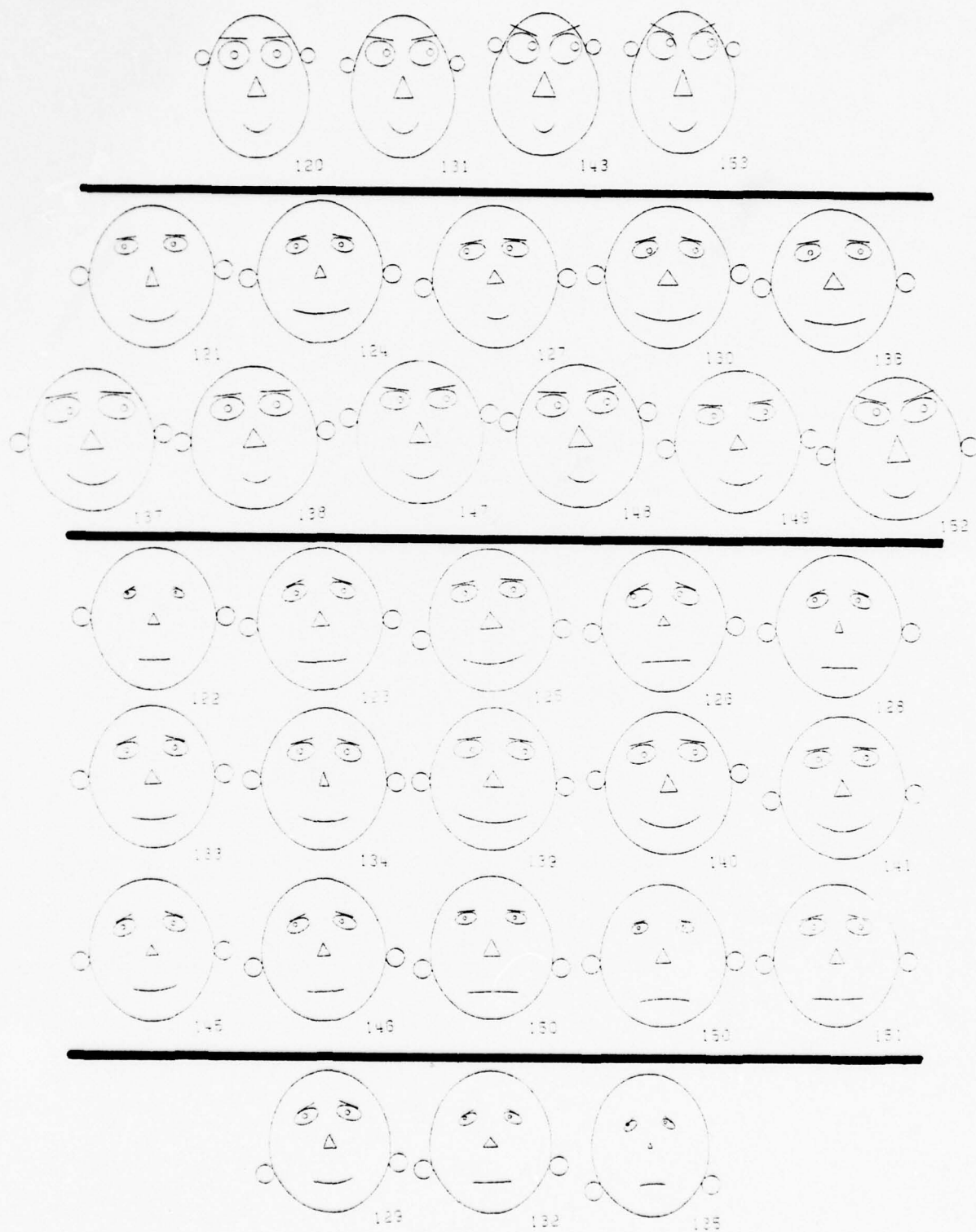


Figure 11

CLUSTERS DETERMINED BY MIKCA

used as a standard against which the judges' results were compared. Table 7 shows the results of this experiment. There was considerable agreement among the judges, as indicated by the comparison coefficients. There was also a good deal of similarity between the clusters formed by the judges and those formed by the MIKCA algorithm.

Several comments by the judges indicated the difficulties they encountered. The most prevalent comment was the difficulty in deciding which feature to consider the most important. One judge considered the mouth first in every case while another judge used the slant of the eyes as a more important variable. The judges also indicated that trying to evaluate simultaneously differences in many features was quite difficult. It is interesting to note that the judges' results were quite similar despite the fact that different criteria were employed as they formed the clusters.

The SOF identification numbers have been altered for this report. There were two course segments which erroneously reported the same SOF number (see face 150). As one looks at the faces with the discriminant space in mind, it is much easier to form a clustering which is similar to the MIKCA solution. One would be aware, for example, that the position of the pupils is critical in that it can diminish the effect of the smile and impact heavily on the horizontal dimension. This effect can be seen by referring to faces 139 and 140; they are included in a group whose smiles are

TABLE 7
COMPARISON COEFFICIENTS FOR ALL PAIRS OF JUDGES

		JUDGES					
		1	2	3	4	5	6
1	1.0	.69	.73	.82	.73	.78	
2		1.0	.90	.80	.77	.68	
3			1.0	.65	.81	.67	
4				1.0	.68	.73	
5					1.0	.69	
6							1.0

COMPARISON COEFFICIENTS
BETWEEN EACH JUDGE AND MIKCA

<u>JUDGE</u>	<u>COMPARISON COEFFICIENT</u>
1	.59
2	.58
3	.68
4	.81
5	.76
6	.73

SIMULTANEOUS COMPARISONS OF
MULTIPLE JUDGES

<u>NUMBER OF JUDGES</u>	<u>COMPARISON COEFFICIENT</u>
3	.52
4	.44
5	.38

not as large as their own because the position of their pupils (positive to the right) has diminished the impact of the curvature of the mouth.

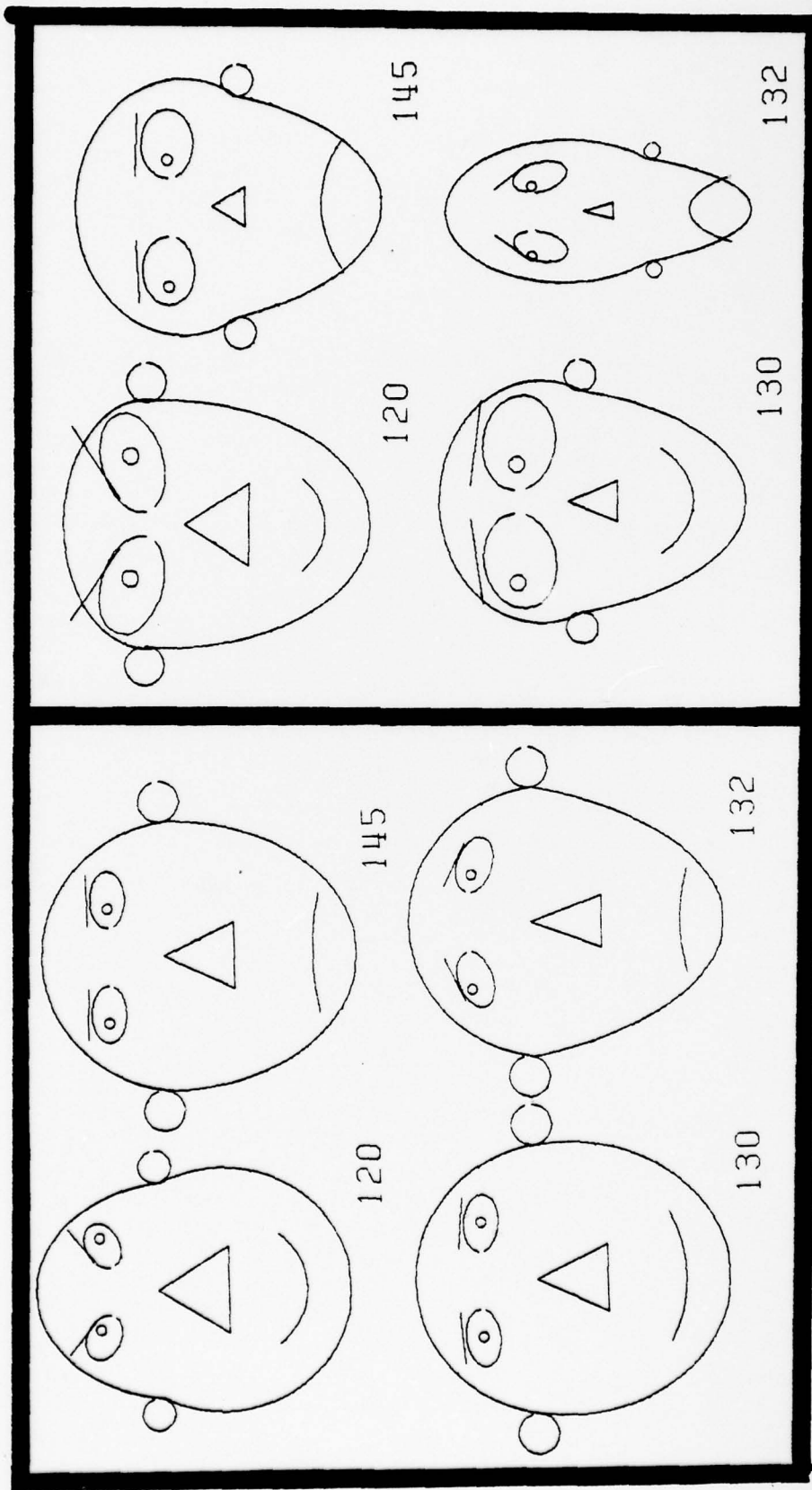
Another example of the interaction of variables is seen by referring to face 152. Most judges would quickly include this face with the group of four at the top of figure 11. A subtle difference, however, is the ear position. Reference to the discriminant function coefficients will indicate a negative which has offset the slant of the eyes in the second dimension.

Knowledge of the discriminant functions helps alleviate the confusion which sets in when attempting to cluster. It is especially true in this case where so little difference exists between the majority of the faces in the middle groups.

Difficulty in evaluating all 13 features simultaneously was a problem. As an alternative to this set of faces, two other sets were produced, one with only six variables features and the other with eight. Figure 12 contains samples from these sets, 12a the six variable set and 12b the eight variable. Of course, not all of the data is represented in this manner. Only those variables which loaded heavily in the discriminant analysis were used, and the features controlled by those variables are the ones considered to convey the most information. Table 6 gives the complete feature-variable combinations used in the construction of all sets of faces. The data used in constructing the set of 33 faces is found in Appendix B.

SIX VARIABLES

LIGHT VARIABLES



12a

12b

Figure 12

D. PROBLEMS ENCOUNTERED

The last section addressed difficulties faced by the judges because it was impossible for them to be aware of the information contained in the discriminant analysis. Of course, it would be pointed out that there is little reason to use this particular MIKCA solution as the standard, but it does serve as an objective standard, as it was desired to compare the machine results with the human results. This section addresses problems of a more mechanical nature.

Exploratory work with the faces uncovered quite a number of relationships between the features. The existence of geometric dependencies (not discriminant-type effects) between features caused difficulties in clearly displaying the variables. Two notable examples are mentioned here.

The length of the mouth is quite dependent upon its curvature. The projection on the horizontal axis (no relation to discriminant axis) has half-length

$$a_m = X9(h/|X8|)$$

where X8 is the mouth curvature. The variables which control these features are thus automatically forced into this dependent status regardless of their true relationship.

The other example concerns the ellipses forming the facial boundary and the angle theta. The upper ellipse is

drawn through the points P', U, and P; the lower through P', L, and P (see figure 10). Two faces with identical values for the ellipses might have quite different appearing facial boundaries due to the dependence on theta for the points P' and P. This is another example of forced dependence.

In order for the width and height of the face to meet a specified constant, the program "normalizes" both horizontal and vertical axes. This normalization eliminates the effects of X1 and X3, and it adjusts all of the features during the process. It is believed to be this normalization process which causes faces which are growing wider and wider to suddenly revert to one-half the widest width when the width exceeds a threshold value. A similar phenomenon is experienced in the height variable. This half-size adjustment may be seen in figure 12b. Face 132 has been changed by a disproportionate amount due to the normalization process. It is of interest to point out that the face-width feature was being held constant during the construction of that set of faces.

Yet another hidden dependency is that of the nose length on the eye height. The eyes are located at height

$$y_e = h[X_{10} + (1 - X_{10})X_6]$$

where X6 controls the length of the nose.

These and other subtle dependencies mislead the investigator if he is not aware of their existence. These problems

reduce the ability of the faces to effectively display the full 20 dimensions. Unfortunately, these points are not explicit in the original document [11] and their discovery was an 11-th hour surprise. It was not possible to adjust for them or to uncover all such relationships at this writing. Appendix C gives a complete listing of the formulas used in the construction of the faces.

VII. COMPARISON COEFFICIENT

A. BACKGROUND AND ALTERNATIVES

Repeated use of the comparison coefficient has been made in this study. The present chapter is devoted to an explanation of this measure of association. The method should be flexible enough to handle multiple comparisons simultaneously, thus enabling one to measure the overall agreement of several judges.

It was decided the best way to display the agreement of two judges was through the use of a contingency table. Table 8 is an example of one to be used in the discussion.

		Judge X		
		A	B	C
Judge Y	A	5	0	1
	B	1	3	3

Table 8

The contingency table indicates the agreement of the two judges. The purpose of this chapter is to find a measure which evaluates how close this agreement is. Note that judge X categorized the observations into three clusters with 6, 3, and 4 elements, respectively. Of the 13 observations, judge Y placed six in one group and seven in another.

The labeling of the clusters is arbitrary. The upper left entry in the table indicates that five of the objects in judge X's category A matched with five of judge Y's category A. The entire table is interpreted in this manner. Notice that if one chooses to call this entry of five as representing agreement, then the entry of 1 below it and the 1 in the top right corner must represent some of the observations on which the judges disagreed.

The contingency table idea is easily generalized to higher dimensions (more than two judges). In three dimensions, a box (or cube) would represent the table, with elements internal to the box measuring agreement between three judges.

One method for measuring the degree of agreement is to find the largest combination of entires such that only one per row and one per column are chosen. This task becomes very difficult as the number of clusters increases, but it can be solved through the use of linear programming techniques. It is a constrained optimization with a linear objective function and is an application of the "assignment problem." Unfortunately, when generalizing to higher dimensions, the L.P. loses its unimodularity attribute and the number of constraints and variables in the problem becomes prohibitively large.

The Chi-square contingency statistic was considered inappropriate because, when using the smaller sample sizes, more than 20 percent of the cells have expected frequencies

of less than five (see ref. 14). Even when using the 190 element sample, there were frequent occasions when this same difficulty persisted. The Chi-square statistic was not used since it could not have been applied consistently throughout the analysis.

Professor James Hartman provided an idea that led to the method finally put into use.

B. THE TECHNIQUE

The idea was to sum the squares of the entries in the contingency table and then "normalize" this quantity. Summing squares offers an excellent method for measuring the degree of association, however the following example illustrates the need for some sort of adjustment factor.

10	0
0	10

9a

19	0
0	1

9b

TABLE 9

Both tables represent perfect agreement on 20 observations, however table 9a has a sum of squares equal to 200 and 9b has a value of 362. It is desired to indicate both of these examples as perfect agreement with one being no better than the other. Hence, it became necessary to determine the "best possible" sum of squares in every given

situation. A computer program was written for this purpose and is included in Appendix F. The statistic which is used as a comparison coefficient is a number between zero and one, formed as the ratio of the actual sum of squares to the "best possible" sum of squares. The best possible sum of squares is a computed sum using a minimax approach and is based on the number of judges, number of clusters by each judge, and the number of observations within each cluster. The minimax procedure does not need to know which observations make up a cluster, only how many observations. An example showing the computation of the comparison coefficient is given in Appendix E.

This method for measuring the degree of agreement provides the analyst a standard scale upon which to compare coefficients based on solutions involving varying numbers of clusters and cluster memberships, as well as varying numbers of judges.

In order to provide some sensitivity for the significance of this measure, several cluster solutions were formed wholly at random and compared to results produced by MIKCA and the judges. In every case, the values of the comparison coefficients were less than 0.1.

VIII. SUMMARY AND CONCLUSIONS

This research has been largely exploratory. A path has been paved for others to follow in examining the SOF data. The theory of cluster analysis and its relationship to discriminant theory have been carefully examined with emphasis on two widely divergent techniques. In the analysis of the data, attempts have been made to identify the underlying structure of which the clusters are a consequence. This chapter is devoted to separate discussions of the cluster methodology explored and the interpretation of the SOF data.

Although the development of methodology phase of the research was carried to completion in a general sense, a number of problems were encountered along the way. Many of these problems are deserving of deeper treatment and are discussed below.

The data transformation was the best of the three considered. It produced the smallest test statistic for homogeneity of covariances, but the value itself was not in the acceptable range, based on normal theory. It should be possible to improve the choice.

The modifications of MIKCA to allow for weighting of the input vectors has been effected and well tested. It is an important added capability for this program.

The use of discriminant analysis to discover the important variables affecting the clustering is, no doubt,

not new. It needs some refinement, however, because it is not clear how one should rate the importance of variables supporting the first dimension to those supporting the second (or any other dimension). Such a set of priorities could be most useful.

The idea of using the important variables (and their signs) of the discriminant functions in the problem of assigning sets of variables to sets of features is believed to be new. It may have great potential in providing a way for the Chernoff face technique to replace the more expensive technique based on computer iteration.

The present attempt to work with the faces was disappointing. This is due largely to the fact that certain restrictions, truncations, and discontinuities in the movement of the features were not well documented in our sources. Their discovery came as a surprise late in the research and it was not feasible to go back and readjust. Such readjustment is clearly called for and would require a substantial effort in the future studies.

The coefficient of comparison was a new idea and there was insufficient time to explore its properties. What is needed is more investigation in order to interpret its various values (or another measure whose values are interpretable). The comparison measure is also useful in the problem of assigning variables to features when working with faces. The goal is to choose assignments having the property

that the judges are in good agreement when forming the clusters.

The study of the SOF data which was made while developing the cluster methods produced results about student evaluations of courses (and instructors). The results are discussed below.

A principal components analysis of the data swarm of mean vectors showed it to be essentially one dimensional and having the direction of the main diagonal of 13-space. The interpretation of this is that all 13 items are equally important in the students' perception of rating the course and its instructor. On the other hand, this same effect would be produced by careless, perfunctory completion of the forms by many students.

The partitioning of the data into three or four clusters by MIKCA is more or less successful. The clusters are not sharply separated (there are no great voids between them). Study should be made to see how much the density of the data diminishes near the boundaries of the partitions.

Although the main data swarm is essentially one dimensional, it appears useful to use two dimensions to describe the individual partitions after clustering. In doing this, variable 12 (overall rating of the instructor) emerged as most important in the first dimension and variables 3, 8, and 11 giving support in the second. Only one discriminant study is reported here, although several were performed.

Variable 12 appears to have permanence while the other variables often shift in importance.

The cluster profiles which track the cluster centroids over the 13 variables provide a set of (almost) horizontal lines. This supports further the one dimensional interpretation of the data swarm. The result is not sensitive to whether or not the data are standardized.

The results of applying the modified MIKCA did not vary greatly from the results of applying the original MIKCA. Hence the number and composition of the clusters is not disturbed much by the variability in class size.

The relative position of the clusters is strongly and inversely related to class size. The courses that receive uniformly high ratings are associated with the small class sizes and the courses receiving uniformly poor ratings are associated with the large class sizes.

All judges reported use of a hierarchical approach to separate the faces into clusters. Most judges first separated the faces into two groups according to the curvature of the mouth (smile or frown). There was little agreement about which features were important in further subdividing the two main groups, hence some disparity resulted in their final cluster solutions.

The MIKCA procedure is a sophisticated approach to cluster analysis; its results are based on sound statistical theory. The modified version of that computer program is considered

particularly well suited to the SOF data or any other data set possessing the same predetermined class structure. The impact of class size on cluster membership has been emphasized. This important issue may indicate the smaller classes receive artificially inflated SOF scores. Consideration to this fact surely must be given by those who use these scores as a means for evaluating teacher performance.

APPENDIX A

Test for the Equality of Dispersion Matrices of k Groups

Given a sample of k groups and m variables with group dispersion matrices, S_i , ($i = 1, \dots, k$) pooled within-groups dispersion matrix S_w , and total sample observations $N = \sum_{g=1}^k N_g$, Box shows that the hypothesis

$$H_1: \Sigma_1 = \Sigma_2 = \dots = \Sigma_k$$

may be tested by an F statistic developed from

$$A = \ln[|S_w|] \cdot [N - k] - \sum_{i=1}^k (N_i - 1) \cdot \ln(|S_i|)$$

$$B = \frac{\left[\sum_{i=1}^k \frac{1}{1 + N_i} - \frac{1}{N - k} \right] \cdot (2m^2 + 3m - 1)}{6(k - 1)(m + 1)}$$

$$C = \frac{\left[\sum_{i=1}^k \frac{1}{(1 + N_i)^2} - \frac{1}{(N - k)^2} \right] \cdot (m - 1) \cdot (m + 2)}{6(k - 1)}$$

$$D = \frac{(k - 1) \cdot m \cdot (m + 1)}{2}$$

$$E = \frac{D + 2}{\text{abs } B^2 - C}$$

If $B^2 > C$, the test statistic is

$$\left(\frac{E}{D}\right) \left[\frac{A(1 - B + 2/E)}{E - A(1 - B + 2/E)} \right] \sim F_{E,D}$$

If $C > B^2$, the test statistic is

$$\left(\frac{A}{D}\right) (1 - B - D/E) \sim F_{E,D}$$

AD-A068 544

NAVAL POSTGRADUATE SCHOOL MONTEREY CALIF
DEVELOPMENT OF CLUSTER ANALYSIS METHODS SUITABLE FOR STUDENT OP--ETC(U)
MAR 79 J W AIKEN

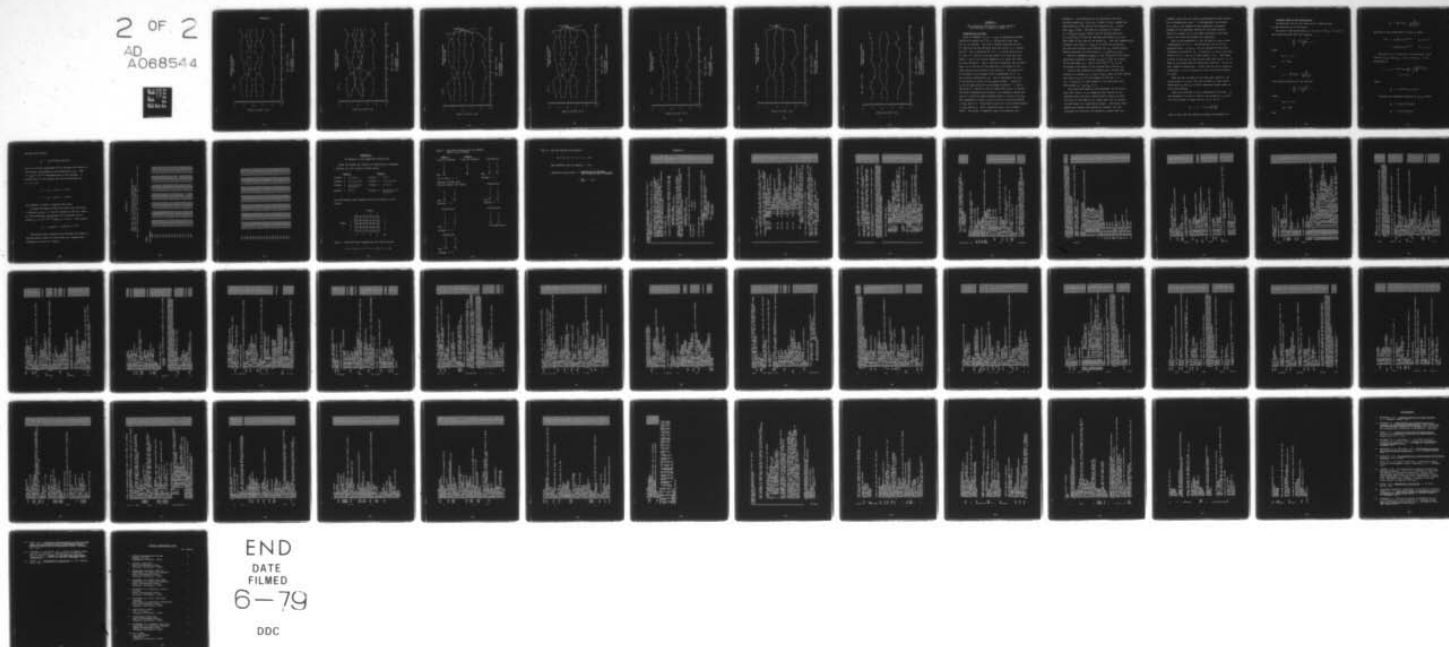
F/G 12/1

UNCLASSIFIED

NL

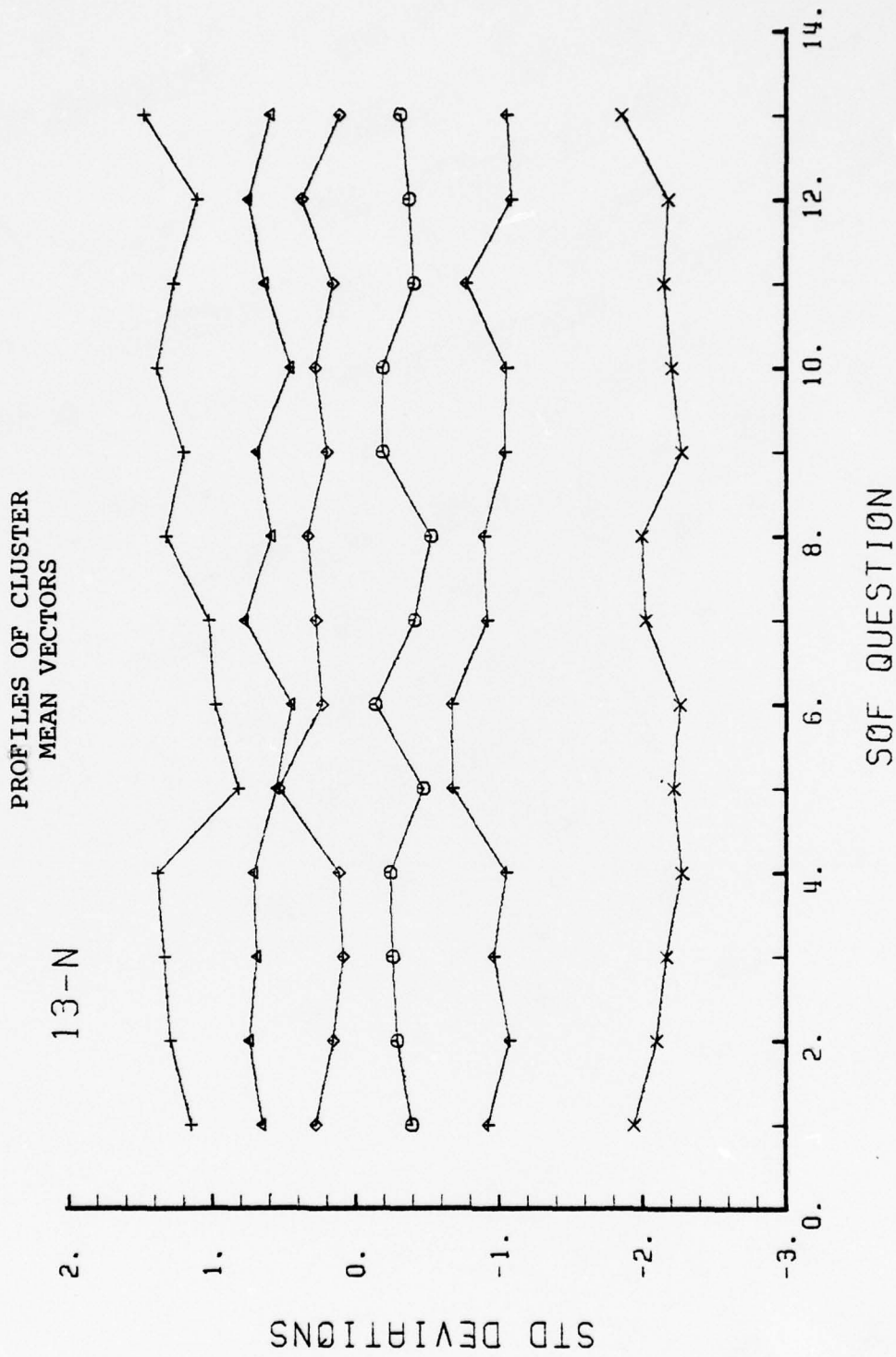
2 OF 2

AD
A068544

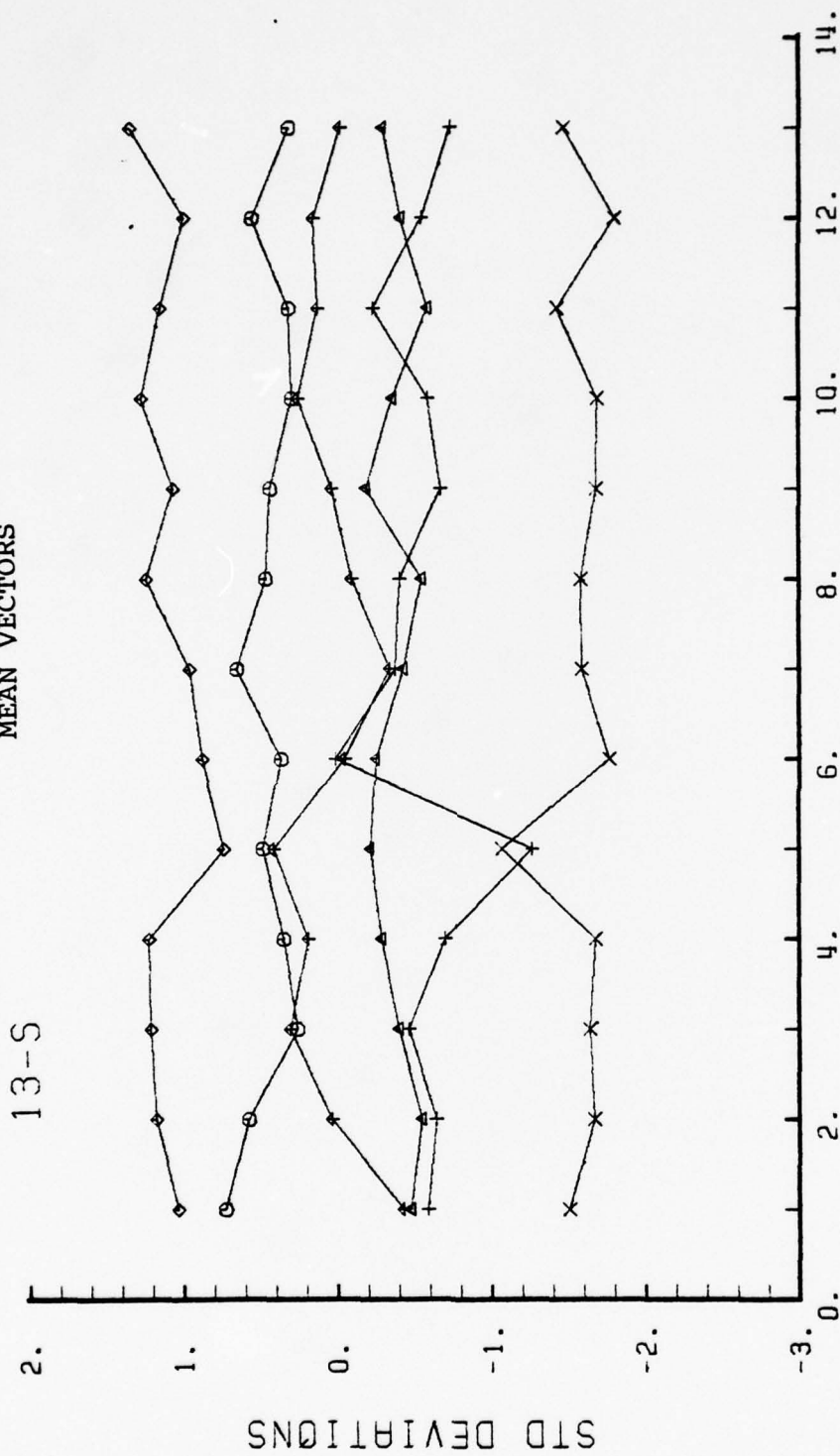


END
DATE
FILMED
6-79
DDC

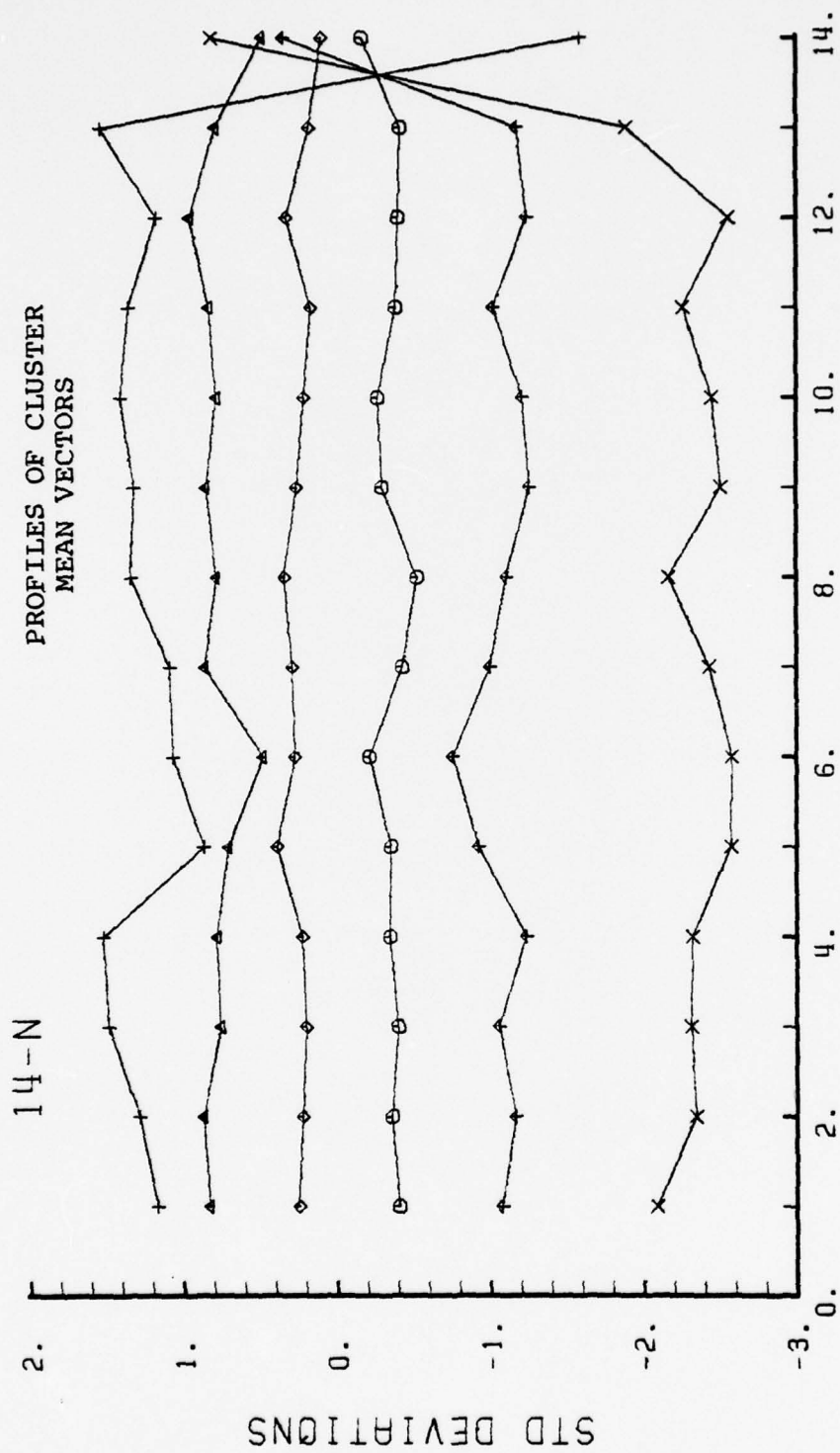
APPENDIX B



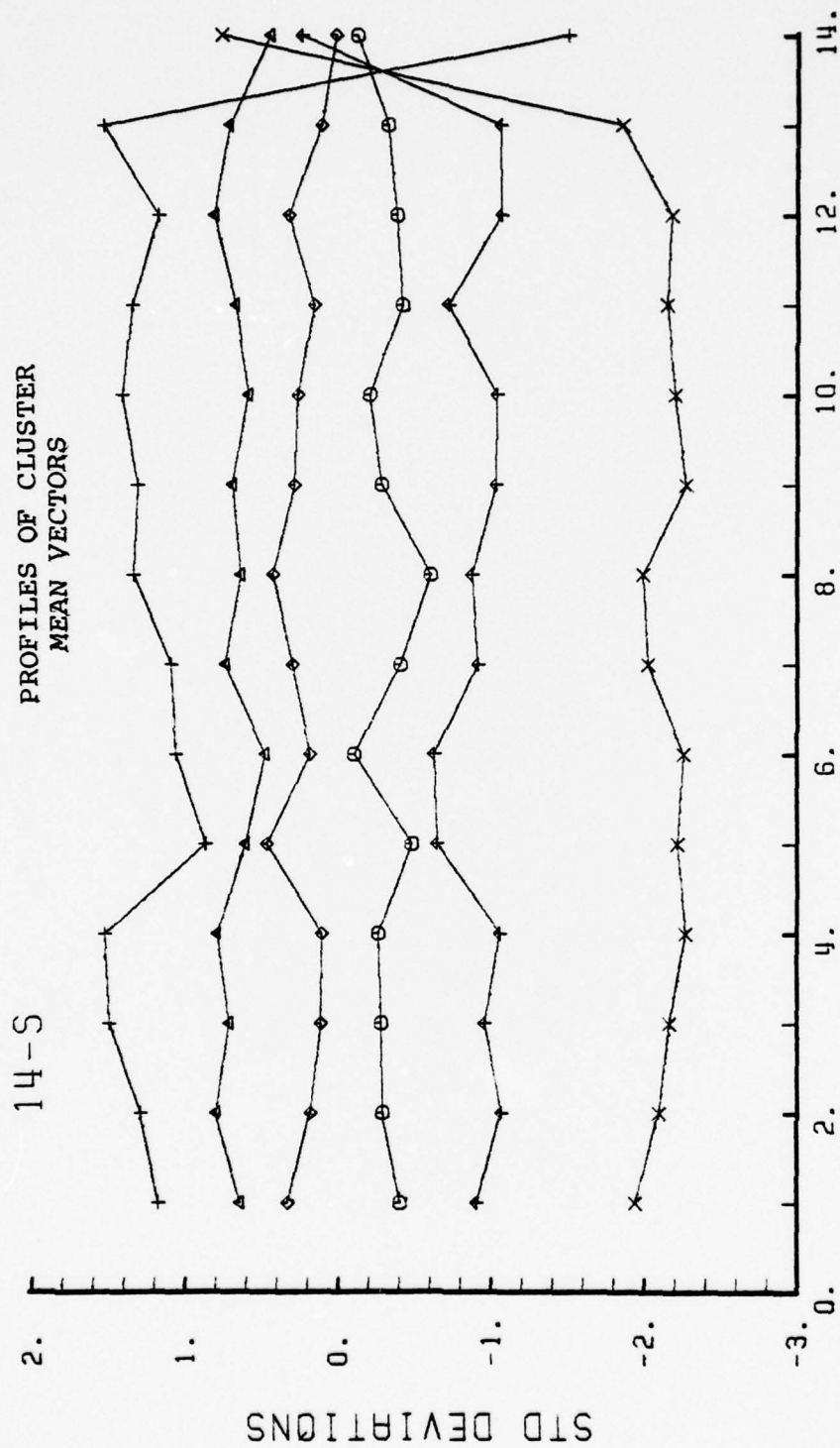
PROFILES OF CLUSTER MEAN VECTORS



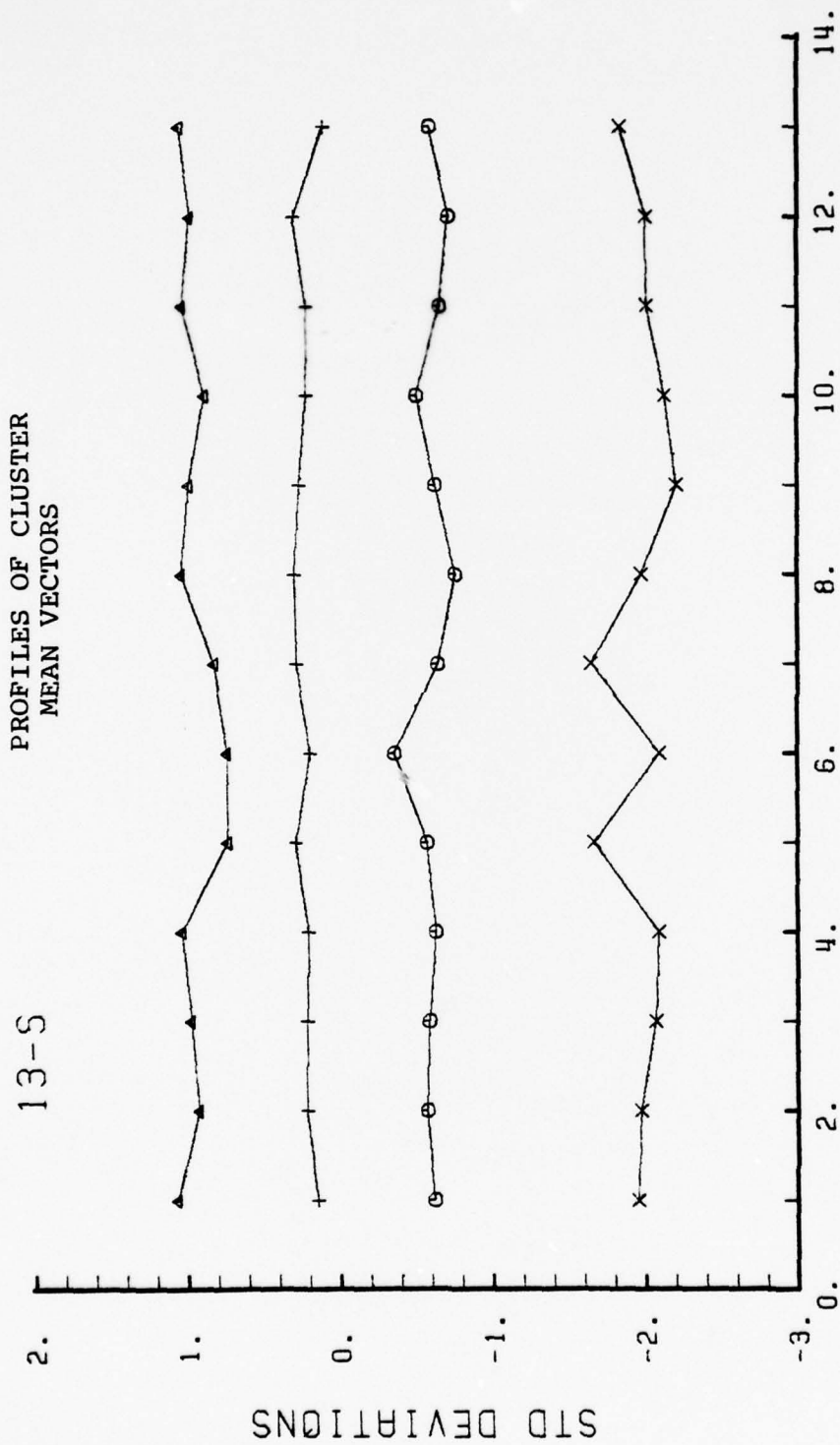
SIX CLUSTERS
13 VARIABLES, STANDARDIZED



SIX CLUSTERS
14 VARIABLES, NOT STANDARDIZED

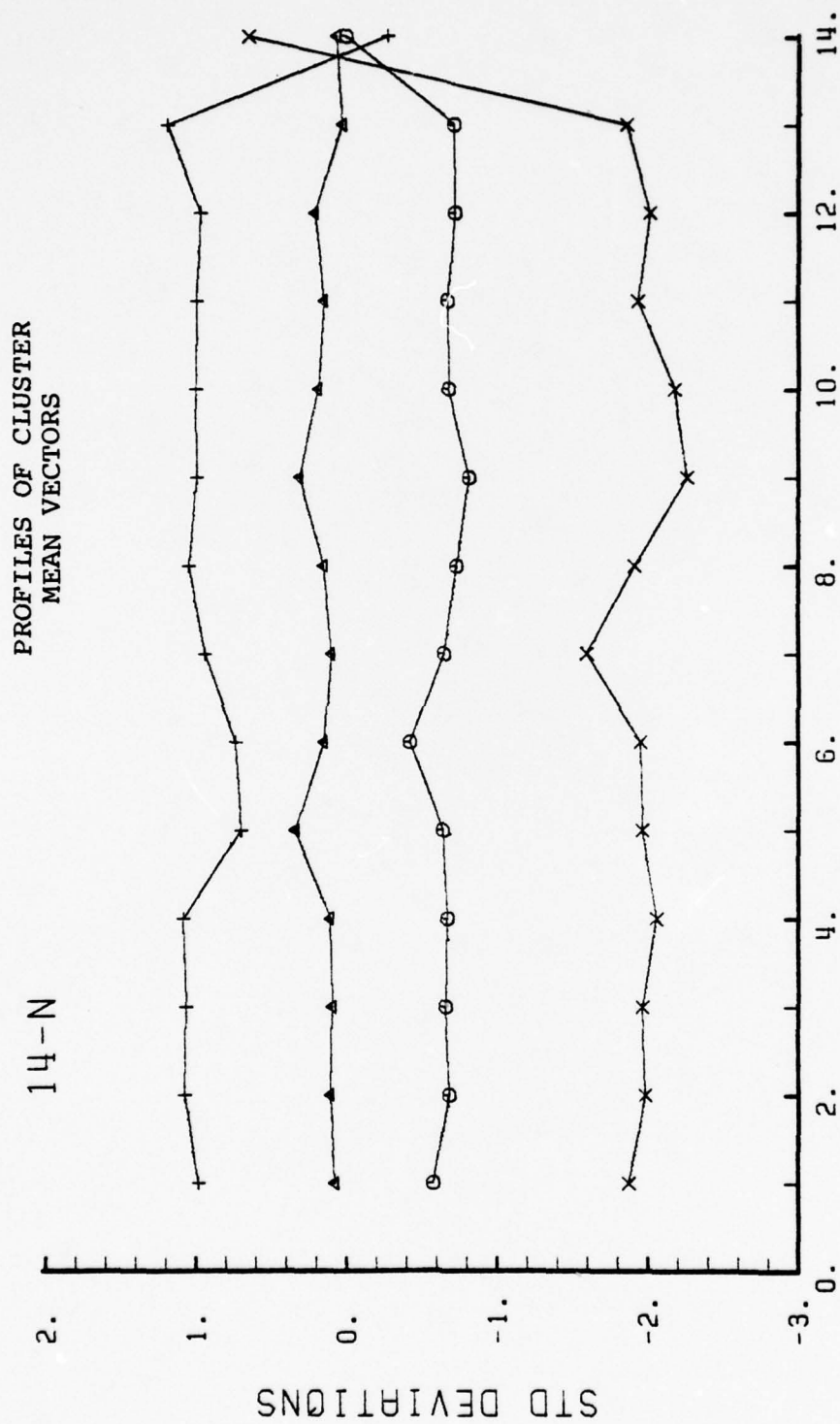


SOF QUESTION
SIX CLUSTERS
14 VARIABLES, STANDARDIZED



SOF QUESTION

FOUR CLUSTERS
13 VARIABLES, STANDARDIZED

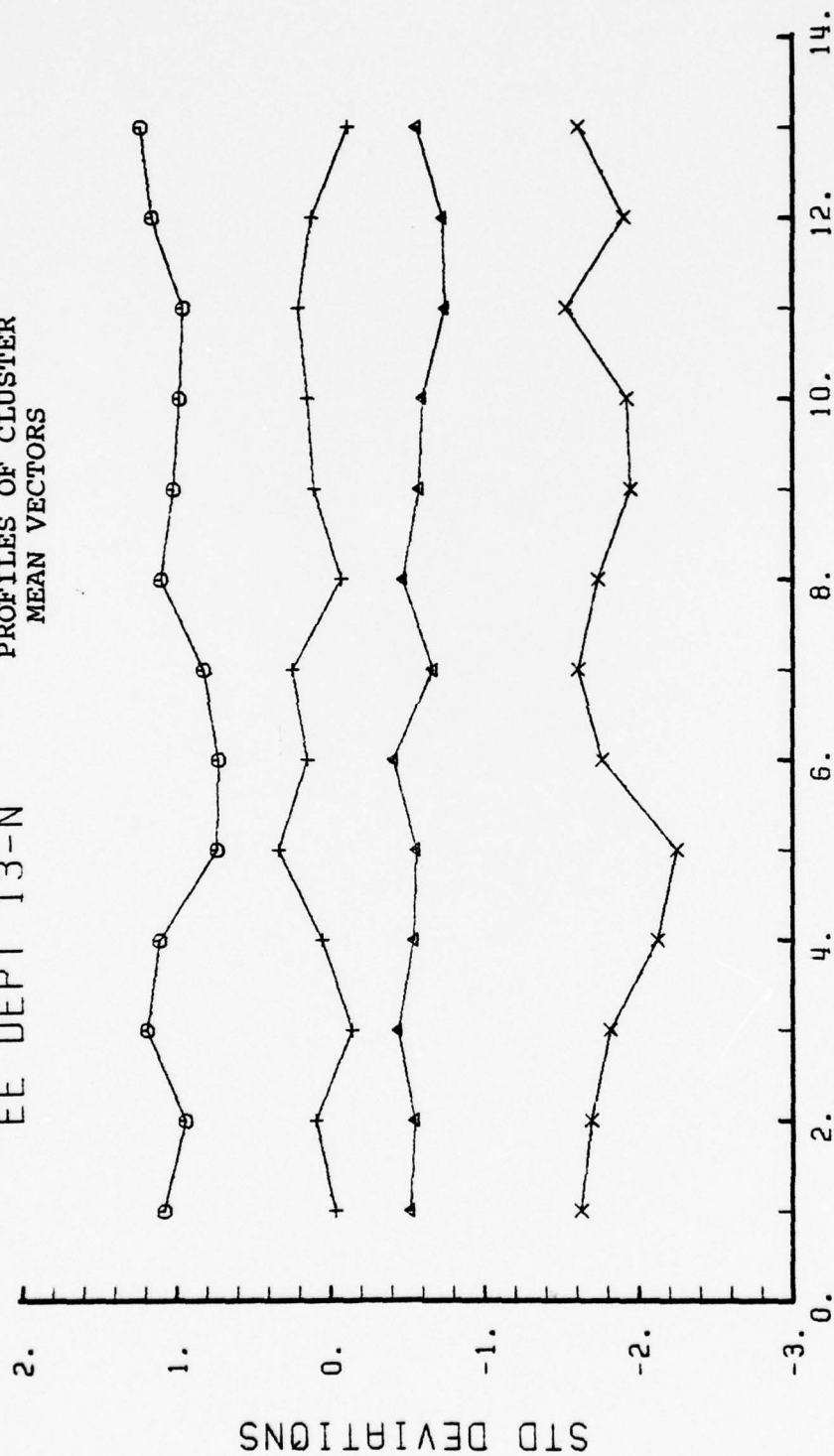


SOF QUESTION

FOUR CLUSTERS
14 VARIABLES, NOT STANDARDIZED

EE DEPT 13-N

PROFILES OF CLUSTER
MEAN VECTORS



SOF QUESTION

FOUR CLUSTERS
13 VARIABLES, NOT STANDARDIZED
DATA FROM ELECTRICAL ENGINEERING DEPARTMENT

APPENDIX C

The following information is taken directly from Chernoff's technical report [9]

Construction of Faces

Given 18 numbers $(x_1, x_2, \dots, x_{18})$ in appropriate ranges (which will usually be 0 to 1), we define a face (see Fig. 3) as follows. Let H be a nominal distance and let $h^* = \frac{1}{2}(1+x_1)H$ be the distance from the origin to a "corner" point P . As x_1 varies from 0 to 1, h^* varies from $H/2$ to H . Let $\theta^* = (2x_2-1)\pi/4$ be the angle of OP with the horizontal. Let P' be a point symmetric to P about the vertical axis through O . Let $h = \frac{1}{2}(1+x_3)H$ represent the distance from O to U the top of the head and L the bottom of the head, both on the vertical line through O . The upper part of the head is an ellipse which is determined by P' , U , and P and an eccentricity x_4 . Let x_4 represent the ratio of the width to height of the upper ellipse. Similarly, x_5 is the same ratio for the ellipse through P' , L , and P . The nose is a vertical line of length $2hx_6$ with O as center. The mouth intersects the vertical line extended through the nose at a point P_m whose distance below O is $h[x_7+(1-x_7)x_6]$. This represents a point x_7 part of the way from the bottom of the nose to U . The mouth is part of a circle whose center is h/x_8 above P_m . Thus a positive value of x_8 yields a smile. The mouth is symmetric about the vertical axis

through 0. Its projection on the horizontal axis has the half-length $a_m = x_9(h/|x_8|)$ unless $(h/|x_8|)$ exceeds the half-width w_m of the face at the height of P_m . In that case $x_9 w_m$ is used. The eyes are located at a height $y_e = h[x_{10} + (1-x_{10})x_6]$ above 0 and at centers which are $x_e = w_e(1+2x_{11})/4$ from the vertical axis where w_e is the half-width of the face at the height y_e . They are symmetrically slanted at an angle $\theta = (2x_{12}-1)\pi/5$ with the horizontal. The eyes are ellipses with eccentricity x_{13} (height/length before slanting) and half-length $L_e = x_{14}\min(x_e, w_e - x_e)$.

The only asymmetry appears in the location of the pupils which move together an amount $r_e(2x_{15}-1)$ from the center of the eye where $r_e = (\cos^2\theta + \sin^2\theta/x_{13}^2)^{-1/2}L_e$ is the horizontal half-length of the slanted eye at height y_e .

Finally the eyebrows are symmetrically located with centers at a height $y_b = 2(x_{16}+.3)L_e x_{13}$ above the eye centers and slant $2(x_{17}-1)\pi/5$ with respect to the eye, i.e., $\theta^{**} = \theta + (2x_{17}-1)\pi/5$ with respect to the horizontal and half-length $L_b = r_e(2x_{18}+1)/2$.

One final step taken by the programmer and which has been left intact, is to normalize both horizontal and vertical axes, each by a multiplicative factor, so that the width of the head at its widest part and its height are both equal to a specified constant. This step, which essentially removes two degrees of freedom, was left unaltered for intuitive and aesthetic reasons that are

somewhat vague and may require reconsideration when dealing with 18-dimensional data. In the meantime, the effects of x_1 and x_3 are almost but not completely eliminated because of the secondary effects of the normalization, which will adjust all of the other features at the same time as the width and height are normalized.

Most of the parameters x_i are adjusted to range within a subinterval of $(0,1)$. The exceptions are two of the eccentricities, x_4 and x_5 , and the parameter controlling curvature of the mouth, x_8 . Ordinarily, x_4 and x_5 are kept within $1/2$ to 2 , and x_8 is kept within $(-5,5)$. The eccentricity of the eye x_{13} has usually been kept within $(.4,.8)$. Some of the ranges must be controlled carefully. We do not want negative length eyes. Others need not be so carefully controlled. It is no calamity to have eyes extend beyond the face.

When the two ellipses of the head meet smoothly, the corner point P is lost, and the variable x_2 loses effect. Restricting x_4 and x_5 to widely separated ranges seems to avoid this problem.

Data are converted to the x parameters as follows. If the variable Z is used to control the parameter x_i , which is to be allowed to range from a_i to b_i , we let

$$x_i = a_i + (b_i - a_i) \left| \frac{Z - m}{M - m} \right|$$

where m and M are the observed minimum and maximum of Z .

Formulae Used on the Construction

We describe a few of the less trivial formulae used in the construction of the faces.

The point P has coordinates $x_o = h^* \cos \theta^*$ and $y_o = h^* \sin \theta^*$. The ellipse through PUP' has equation

$$\frac{x^2}{a_u^2} + \frac{(y - c_u)^2}{b_u^2} = 1$$

where

$$b_u = h - c_u ,$$

$$a_u = x_4 b_u$$

and

$$c_u = \frac{1}{2}[(h+y_o) - \frac{x_o^2}{x_4^2(h-y_o)}] .$$

The ellipse through PLP' has equation

$$\frac{x^2}{a_L^2} + \frac{(y - c_L)^2}{b_L^2} = 1$$

where

$$b_L = h + c_L ,$$

$$a_L = x_5 b_L$$

and

$$c_L = \frac{1}{2} [(-h+y_o) - \frac{x_o^2}{x_5^2 (-h-y_o)}]$$

The head is then described by $(\pm x(y), y)$ where

$$\begin{aligned} x(y) &= x_4 [b_u^2 - (y - c_u)^2]^{1/2} & y_o \leq y \leq h \\ &= x_5 [b_L^2 - (y - c_L)^2]^{1/2} & -h \leq y \leq y_o \end{aligned}$$

The mouth is a circular arc with curvature $|x_8/h|$ through $(0, y_m)$ where $y_m = -h(x_7 + (1-x_7)x_6)$. It is described by

$$\begin{aligned} y &= y_m + (\text{sgn } x_8) \left[\frac{h}{|x_8|} - \sqrt{\left(\frac{h}{x_8}\right)^2 - x^2} \right], \\ 0 &\leq x \leq a_m \end{aligned}$$

where

$$a_m = x_9 \min[x(y_m), h/|x_8|] .$$

The eyes are nominally centered at (x_e, y_e) where

$$y_e = h[x_{10} + (1-x_{10})x_6]$$

$$x_e = x(y_e) [1 + 2x_{11}] / 4$$

and have half-length

$$L_e = x_{14} \min[x_e, x(y_e) - x_e] .$$

Let (u,v) be the coordinates of an ellipse with center at the origin, half-length L_e and eccentricity x_{13} . Then $v = x_{13}(L^2 - u^2)^{1/2}$ describes part of the ellipse. A similar part of the slanted eye can be described for $0 \leq u \leq L$ by

$$x = x_e + u \cos \theta - v \sin \theta$$

$$y = y_e + u \sin \theta + v \cos \theta$$

and symmetry is used to complete both eyes.

To place the pupils within the eyes, both are moved a distance $r_e(2x_{15} - 1)$ from the center of the eye, where r_e , the horizontal half-length of the slanted eye at height y_e , is $(u^2 + v^2)^{1/2}$ when $v/u = \tan \theta$. This yields

$$r_e = L_e (\cos^2 \theta + x_{13}^{-2} \sin^2 \theta)^{-1/2}$$

The program then normalizes all heights and widths by multiplicative factor k/h and $k/\max x(y)$ respectively. Currently k is set at 2 inches.

These are the 33 observations (transformed data) from the Electrical Engineering Department. There are two rows of data per SOf number. The first 13 entries are SOf item scores, the 14-th is class size.

SOF NUMBER

109

[illegible]

0.010202010.10.11132110.11020111000110103221.1

4 1408 1509 1551 1552 1553 1554 1555 1556 1557 1558 1559 1560 1561 1562 1563 1564 1565 1566 1567 1568 1569 1570 1571 1572 1573 1574 1575 1576 1577 1578 1579 1580 1581 1582 1583 1584 1585 1586 1587 1588 1589 1590 1591 1592 1593 1594 1595 1596 1597 1598 1599 1600 1601 1602 1603 1604 1605 1606 1607 1608 1609 1610 1611 1612 1613 1614 1615 1616 1617 1618 1619 1620 1621 1622 1623 1624 1625 1626 1627 1628 1629 1630 1631 1632 1633 1634 1635 1636 1637 1638 1639 1640 1641 1642 1643 1644 1645 1646 1647 1648 1649 1650 1651 1652 1653 1654 1655 1656 1657 1658 1659 1660 1661 1662 1663 1664 1665 1666 1667 1668 1669 1670 1671 1672 1673 1674 1675 1676 1677 1678 1679 1680 1681 1682 1683 1684 1685 1686 1687 1688 1689 1690 1691 1692 1693 1694 1695 1696 1697 1698 1699 1700 1701 1702 1703 1704 1705 1706 1707 1708 1709 1710 1711 1712 1713 1714 1715 1716 1717 1718 1719 1720 1721 1722 1723 1724 1725 1726 1727 1728 1729 1730 1731 1732 1733 1734 1735 1736 1737 1738 1739 1740 1741 1742 1743 1744 1745 1746 1747 1748 1749 1750 1751 1752 1753 1754 1755 1756 1757 1758 1759 1760 1761 1762 1763 1764 1765 1766 1767 1768 1769 1770 1771 1772 1773 1774 1775 1776 1777 1778 1779 1780 1781 1782 1783 1784 1785 1786 1787 1788 1789 1790 1791 1792 1793 1794 1795 1796 1797 1798 1799 1800 1801 1802 1803 1804 1805 1806 1807 1808 1809 1810 1811 1812 1813 1814 1815 1816 1817 1818 1819 1820 1821 1822 1823 1824 1825 1826 1827 1828 1829 1830 1831 1832 1833 1834 1835 1836 1837 1838 1839 1840 1841 1842 1843 1844 1845 1846 1847 1848 1849 1850 1851 1852 1853 1854 1855 1856 1857 1858 1859 1860 1861 1862 1863 1864 1865 1866 1867 1868 1869 1870 1871 1872 1873 1874 1875 1876 1877 1878 1879 1880 1881 1882 1883 1884 1885 1886 1887 1888 1889 1890 1891 1892 1893 1894 1895 1896 1897 1898 1899 1900 1901 1902 1903 1904 1905 1906 1907 1908 1909 1910 1911 1912 1913 1914 1915 1916 1917 1918 1919 1920 1921 1922 1923 1924 1925 1926 1927 1928 1929 1930 1931 1932 1933 1934 1935 1936 1937 1938 1939 1940 1941 1942 1943 1944 1945 1946 1947 1948 1949 1950 1951 1952 1953 1954 1955 1956 1957 1958 1959 1960 1961 1962 1963 1964 1965 1966 1967 1968 1969 1970 1971 1972 1973 1974 1975 1976 1977 1978 1979 1980 1981 1982 1983 1984 1985 1986 1987 1988 1989 1990 1991 1992 1993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017 2018 2019 2020 2021 2022 2023 2024 2025 2026 2027 2028 2029 2030 2031 2032 2033 2034 2035 2036 2037 2038 2039 2040 2041 2042 2043 2044 2045 2046 2047 2048 2049 2050 2051 2052 2053 2054 2055 2056 2057 2058 2059 2060 2061 2062 2063 2064 2065 2066 2067 2068 2069 2070 2071 2072 2073 2074 2075 2076 2077 2078 2079 2080 2081 2082 2083 2084 2085 2086 2087 2088 2089 2090 2091 2092 2093 2094 2095 2096 2097 2098 2099 2100 2101 2102 2103 2104 2105 2106 2107 2108 2109 2110 2111 2112 2113 2114 2115 2116 2117 2118 2119 2120 2121 2122 2123 2124 2125 2126 2127 2128 2129 2130 2131 2132 2133 2134 2135 2136 2137 2138 2139 2140 2141 2142 2143 2144 2145 2146 2147 2148 2149 2150 2151 2152 2153 2154 2155 2156 2157 2158 2159 2160 2161 2162 2163 2164 2165 2166 2167 2168 2169 2170 2171 2172 2173 2174 2175 2176 2177 2178 2179 2180 2181 2182 2183 2184 2185 2186 2187 2188 2189 2190 2191 2192 2193 2194 2195 2196 2197 2198 2199 2200 2201 2202 2203 2204 2205 2206 2207 2208 2209 2210 2211 2212 2213 2214 2215 2216 2217 2218 2219 2220 2221 2222 2223 2224 2225 2226 2227 2228 2229 2230 2231 2232 2233 2234 2235 2236 2237 2238 2239 2240 2241 2242 2243 2244 2245 2246 2247 2248 2249 2250 2251 2252 2253 2254 2255 2256 2257 2258 2259 2260 2261 2262 2263 2264 2265 2266 2267 2268 2269 2270 2271 2272 2273 2274 2275 2276 2277 2278 2279 2280 2281 2282 2283 2284 2285 2286 2287 2288 2289 2290 2291 2292 2293 2294 2295 2296 2297 2298 2299 2300 2301 2302 2303 2304 2305 2306 2307 2308 2309 2310 2311 2312 2313 2314 2315 2316 2317 2318 2319 2320 2321 2322 2323 2324 2325 2326 2327 2328 2329 2330 2331 2332 2333 2334 2335 2336 2337 2338 2339 2340 2341 2342 2343 2344 2345 2346 2347 2348 2349 2350 2351 2352 2353 2354 2355 2356 2357 2358 2359 2360 2361 2362 2363 2364 2365 2366 23

0.000011120111101237110011212010000000220

0.00.00.10.10.00.00.10.10.22.10.00.00.21.10.10.00.00.10.00.00.33.33.
2799.889.1640.6392.8630.9536.0445.7249.7457.1453.1903.2593.6717.8723.9808.6969.2506.1969.7747.8888.0444.
54425471.6487.4522.3640.5587.2597.1530.3671.2256.9387.1882.0262.2225.5555.
84425471.6487.4522.3640.5587.2597.1530.3671.2256.9387.1882.0262.2225.5555.

[illegible][illegible]

135
136
137
138
139
140
141
143
145
146
147
148
149
150
150
151
152
153

APPENDIX E

An Example of the Comparison Coefficient

Given two judges who cluster 20 observations (numbered 1 through 20) into groups as shown below:

<u>Judge X</u>		<u>Judge Y</u>	
Cluster 1	1,2,3,4,7	cluster 1	5,6,7
Cluster 2	5,6,11,12,13	cluster 2	1,2,3,4,9,15
Cluster 3	8,9,10,14,15, 16,17,18	cluster 3	8,11,12
Cluster 4	19,20	cluster 4	10,13,14,16,17, 18,19,20

The contingency table appears below with marginal (row) totals.

		Judge X				
		1	2	3	4	
Judge Y	1	1	2	0	0	3
	2	4	0	2	0	6
	3	0	2	1	0	3
	4	0	1	5	2	8
		5	5	8	2	20

Step 1: Find the sum of squares of the table entries.

$$1 + 4 + 16 + 4 + 4 + 1 + 1 + 25 + 4 = 60$$

Step 2: Find best possible sum of squares.
(Read in two columns)

<u>Judge X</u>		<u>Judge Y</u>		Subtracting	
# obs in clusters		# obs in clusters			
	3		5	0	0
	6		5	1	0
	3		8	1	0
	<u>8</u>		<u>2</u>	<u>0</u>	<u>2</u>
Max	8		8	Max	1

Min of Max's = 8

Minimax = 1

Subtract Minimax from
the max element and repeat

				Subtracting	
	3	5		0	0
	6	5		0	0
	3	0		1	0
	<u>0</u>	<u>2</u>		<u>0</u>	<u>1</u>
Max	6	5		Max	1
Minimax	= 5			Minimax	= 1

Subtracting				Subtracting	
	3	0		0	0
	1	5		0	0
	3	0		0	0
	<u>0</u>	<u>2</u>		<u>0</u>	<u>0</u>
Max	3	5		Finished Step 2	
Minimax	= 3				

Subtracting	
	0
	1
	3
	<u>0</u>
Max	3
Minimax	= 2

Step 3: Sum the squares of Minimax's

$$54 + 25 + 9 + 4 + 1 + 1 = 104$$

$$\text{Best possible sum of squares} = 104$$

$$\text{Comparison coefficient} = \frac{\text{Actual sum of squares}}{\text{Best Possible Sum of Squares}}$$

$$= \frac{60}{104} = 0.58$$

APPENDIX F

```

///AIKSWKEE JOB (1642,0526,RJ74),*AIKEN SMC 2132',TIME=5
/// EXEC FORTC LG,REGION.GO=200K
///FORT.SYSIN DD *

```

K-MEANS ITERATIVE CLUSTERING PROGRAM BY D.J. MCRAE

THIS PROGRAM ENABLES THE USER TO CLUSTER A DATA MATRIX UP TO SIZE (600 X 20) INTO A MAXIMUM OF 20 CLUSTERS. SEVERAL OPTIONS ARE AVAILABLE WITH RESPECT TO THE METHOD OF COMPUTATION AND CRITERIA TO ACCOMPLISH THE CLUSTERING. THE PROGRAM MAY BE USED ON EITHER THE CP/CMS TERMINALS OR BY THE CARD READER. THE APPROPRIATE METHOD FOR EACH IS DESCRIBED BELOW.

CP/CMS TERMINAL USE

THE PROGRAM SHOULD BE READ INTO CP AS IS, BUT PRECEDED BY THE FOLLOWING CARDS, STARTING IN COLUMN ONE

CP67USERID XXXXG
OFFLINE READ RRREAD FORTRAN

THE DATA MATRIX CAN BE PLACED IN THE COMPUTER EITHER BY USE OF THE OFFLINE READ OR BY TYPING IN THE INFORMATION ON THE TERMINAL. IF OFFLINE READ IS USED THE DATA MUST BE IN THE SPECIFIED FORMAT PRECEDED BY THE FOLLOWING TWO CARDS, STARTING IN COLUMN ONE:

CP67USER ID XXXX;
OFFLINE READ FILE FT04FYYY

THE PROGRAM CAN THEN BE EXECUTED BY FIRST ISSUING THE
CMS COMMAND F RRREAD
FOLLOWED BY

OS/CARD READER USE

```

THE PROGRAM MUST CONTAIN THE FOLLOWING INFORMATION
IN THE FOLLOWING ORDER:
STANDARD GREEN JOB CARD, TIME=(AS DESIRED)
// EXEC FORTCLG, REGION.GO=200K
// FORT.SYSIN DD *
MAIN PROGRAM
/*
//GO.SYSIN DD *
DATA DECK

```

RRR00010
RRR00020
RRR00030
RRR00040
RRR00050
RRR00060
RRR00070
RRR00080
RRR00090
RRR00100
RRR00110
RRR00120
RRR00130
RRR00140
RRR00150
RRR00160
RRR00170
RRR00180
RRR00190
RRR00200
RRR00210
RRR00220
RRR00230
RRR00240
RRR00250
RRR00260
RRR00270
RRR00280
RRR00290
RRR00300
RRR00310
RRR00320
RRR00330
RRR00340
RRR00350
RRR00360
RRR00370
RRR00380
RRR00390
RRR00400
RRR00410
RRR00420
RRR00430
RRR00440
RRR00450

RRR00460
RRR00470
RRR00480
RRR00490
RRR00500
RRR00510
RRR00520
RRR00530
RRR00540
RRR00550
RRR00560
RRR00570
RRR00580
RRR00590
RRR00600
RRR00610
RRR00620
RRR00630
RRR00640
RRR00650
RRR00660
RRR00670
RRR00680
RRR00690
RRR00700
RRR00710
RRR00720
RRR00730
RRR00740
RRR00750
RRR00760
RRR00770
RRR00780
RRR00790
RRR00800
RRR00810
RRR00820
RRR00830
RRR00840
RRR00850
RRR00860
RRR00870
RRR00880
RRR00890
RRR00900
RRR00910
RRR00920
RRR00930

/*

DATA INPUT DECK

THE FIRST CARD OF THE DATA DECK IS THE TITLE CARD.
IT MAY CONTAIN ANY ALPHA NUMERIC TITLE IN COLUMNS 1 - 80.

THE SECOND CARD IS THE PROBLEM CARD. IT CONTAINS
INTEGERS IN THE FIRST 13 COLUMNS IN THE FOLLOWING MANNER:

COL 1-4: NUMBER OF OBSERVATIONS
COL 5-6: NUMBER OF VARIABLES PER OBSERVATION
COL 7-8: NUMBER OF CLUSTERS DESIRED
COL 9: CRITERION TO BE USED IN THE EVALUATION:

1 = TRACE W
2 = DETERMINANT W
3 = LARGEST ROOT OF W(INVERSE)*B
4 = TRACE W(INVERSE)*B

COL 10: STANDARDIZATION PARAMETER:

0 = THE PROGRAM WILL NOT STANDARDIZE DATA
1 = THE PROGRAM WILL STANDARDIZE DATA

COL 11:

1 = DISTANCE PARAMETER
0 = EUCLIDIAN DISTANCE
1 = SCALED EUCLIDIAN DISTANCE
2 = MAHALANOBIS DISTANCE

COL 12:

DATA
0 = DATA IS TO BE READ IN
1 = NO DATA READ IN: PREVIOUS DATA IS
TO BE REANALYZED USING CURRENT PROBLEM
PARAMETER

COL 13:

TIMING
0 = NO OBSERVATIONS CONSIDERED IN
INDIVIDUAL SWITCHES
1 - 8 = AN INCREASING NUMBER OF OBSERVATIONS
9 = ALL OBSERVATIONS USED

NOTE: THE TIMING PARAMETER DETERMINED HOW EXTENSIVE THE
DOUBLE CHECK OF THE CLUSTER SOLUTION IS TO BE.
NORMALLY '9' IS USED UNLESS THE DATA MATRIX IS LARGE

THIS CARD MUST BE IN (14,212,511) FORMAT

THE THIRD CARD IS THE VARIABLE FORMAT CARD. THIS CARD
MUST BE PRESENT UNLESS THE DATA PARAMETER IN THE PROBLEM
CARD IS '1'. THE FORMAT STATEMENT IS THAT TYPE WHICH THE DATA
CARDS ARE IN. FOR EXAMPLE IF THE DATA CARD HAS DATA:

ABCD 1.59
STARTING IN COL 1, THE VARIABLE FORMAT CARD WOULD
CONTAIN THE FOLLOWING IN COL 1:
(1A4,4X,F3.1,4X,F4.2)

CC


```

RRR00940
RRR00950
RRR00960
RRR00970
RRR00980
RRR00990
RRR01000
RRR01010
RRR01020
RRR01030
RRR01040
RRR01050
RRR01060
RRR01070
RRR01080
RRR01090
RRR01100

```

```

THE NEXT CARDS ARE THE DATA CARDS. THEY MAY CONTAIN
DATA IN COLUMNS 1 - 72.

FOR REANALYSIS OF THE SAME DATA, ADD A NEW TITLE CARD
FOLLOWED BY A NEW PROBLEM CARD, WITH A '1' IN COLUMN 11
OF THE NEW PROBLEM CARD.
FOR ANALYSIS OF NEW DATA, SIMPLY PLACE THE NEW DATA
SET, INCLUDING NEW TITLE CARD, PROBLEM CARD, AND FORMAT CARD,
AFTER THE PREVIOUS SET

FOR THE PROGRAM TO EXIT NORMALLY, TWO BLANK CARDS MUST FOLLOW
THE LAST DATA BATCH.

COMMON NOBS, NVAR, NGPS, ICRIT, NOSTAN, IDIST, IFINE, KTIME, IDENT(600),
1 IDATA(600,20), T(20,20), B(20,20), W(20,20), WFCT(20,20), SVCEN(20,20),
2 IDATA(600), NISV(20), VMEAN(20), SD(20), XVEC(20), YVEC(20), BT(20,20),
3 NISVT(20), SVCENT(20,20), IDATA(600), VEC(20,20), EIG(20), NS(600),
4 TOP(20,20), BOT(20), UP(20,20), DOWN(20)

```

```

RRR01120
RRR01130
RRR01140
RRR01150
RRR01160
RRR01170
RRR01180
RRR01190
RRR01200
RRR01210
RRR01220
RRR01230
RRR01240
RRR01250
RRR01260
RRR01270
RRR01280
RRR01290
RRR01300
RRR01310
RRR01320
RRR01330
RRR01340
RRR01350
RRR01360
RRR01370
RRR01380
RRR01390
RRR01400

```

DESCRIPTION OF COMMON AREA:

```

NOBS = NUMBER OF OBSERVATIONS
NVAR = NUMBER OF VARIABLES
NGPS = NUMBER OF CLUSTERS
ICRIT = CRITERION TO BE OPTIMIZED (SEE ABOVE)
NOSTAN = STANDARDIZATION PARAMETER (SEE ABOVE)
IDIST = DISTANCE PARAMETER (SEE ABOVE)
IFINE = ESCAPE PARAMETER: IF IFINE IS SET EQUAL TO '1',
SOMETHING IS WRONG AND THE APPROPRIATE ERROR
MESSAGE IS PRINTED OUT; THE PROGRAM GOES ON TO THE
NEXT PROBLEM

KTIME = TIMING PARAMETER (SEE ABOVE)
IDENT = OBSERVATION IDENTIFICATIONS (A4 FORMAT): READ FROM
EACH DATA CARD

DATA = AREA IN WHICH THE DATA MATRIX IS STORED
T = CROSS-PRODUCTS MATRIX
B = BETWEEN-CLUSTERS MATRIX
W = WITHIN-CLUSTERS MATRIX
WFCT = CHOLESKY FACTOR OF THE WITHIN-CLUSTERS MATRIX
SVCEN = CLUSTER CENTERS (MEANS)
IDATA = CLUSTER IDENTIFICATION FOR EACH OBSERVATION
NISV = CLUSTER SIZES (NUMBER OF OBSERVATIONS IN EACH CLUSTER)
VMEAN = VARIABLE MEANS
SD = VARIABLE STANDARD DEVIATIONS
XVEC = VARIABLE STORAGE
YVEC = TEMPORARY STORAGE

```

```

CCCCCCCCCCCCCCCC

```

```

CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC

```


RRR01410
RRR01420
RRR01430
RRR01440
RRR01450
RRR01460
RRR01470

BT = TEMPORARY STORAGE FOR BETWEEN-CLUSTERS MATRIX
NISVT, SVCENT, IDATAT: TEMPORARY STORAGE SERVING THE
SAME FUNCTIONS AS NISV, SVCEN, AND IDATA
VEC = EIGENVECTORS
EIG = EIGENVALUES

NNBS=0
NNN IS USED 3 TIMES IN NEXT SEVERAL CARDS.

CHECK THIS CARD FOR CORRECT NOBS

NNN=190
READ(5,899)(NS(I),I=1,NNN)

FORMAT(2513)
WRITE(6,888)(NS(I),I=1,NNN)

FORMAT(2513)
DO 200 I=1,NNN

NNBS=NNBS+NS(I)

WRITE(6,887) NNBS

FORMAT(0,0,0,TOTAL,NNBS)

CALL PRELIM(CRIT,NNBS)

IFINE IS ESCAPE PARAMETER

IF ((IFINE.EQ.1)) GO TO 400

IF ((IFINE.EQ.2)) GO TO 500

CALL RANDST (CRIT)

IF ((IFINE.EQ.1)) GO TO 400

CALL KMEANS (CRIT)

IF ((IFINE.EQ.1)) GO TO 400

CALL ISWITCH (CRIT)

GO TO 100

IFINE GOT SET TO 1: IF DATA HAS BEEN READ IN, RESET DATA AND

GO TO NEXT PROBLEM

DO 410 I=1,NOBS

DO 410 J=1,NVARS

IF (NOSTAN.EQ.1) DATA(I,J) = DATA(I,J) * SD(J)

DATA(I,J) = DATA(I,J) + VMEAN(J)

GO TO 100

ALL DONE: WRITE OUT LAST MESSAGE AND EXIT

WRITE(6,900)

FORMAT(10 END OF ANALYSES.)

STOP

END

SUBROUTINE PRELIM(CRIT,NNBS)

THIS SUBROUTINE MAKES THE PRELIMINARY CALCULATIONS.

IT INPUTS THE DATA, CALCULATES THE MEANS AND VARIANCES FOR EACH

VARIABLE, STANDARDIZES (CONVERTS TO Z-SCORES) EACH VARIABLE IF

REQUESTED, AND CALCULATES THE CROSS-PRODUCTS MATRIX

RRR01490
RRR01500
RRR01510
RRR01520
RRR01530
RRR01540
RRR01550
RRR01560
RRR01570
RRR01580
RRR01590
RRR01600
RRR01610
RRR01620
RRR01630
RRR01640
RRR01650
RRR01660
RRR01670
RRR01680
RRR01690

RRR01710
RRR01720
RRR01730
RRR01740
RRR01750
RRR01760

C
C
C
C
C
C

C
C

899

888

200

987

100
C

300

C
C

400

410

C
500
900

C
C
C
C
C

```

COMMON NOBS,NVARS,NGPS,ICRIT,NOSTAN,IDIST,IFINE,KTIME,IDENT(600),
1 DATA(600,20),T(20,20),B(20,20),W(20,20),WECT(20,20),SWCEN(20,20),
2 INDATA(600),NISV(20),VMEAN(20),SD(20),XVEC(20),YVEC(20),BT(20,20),
3 NISVT(20),SVCENT(20,20),IDATAT(600),VEC(20,20),EIG(20),NS(600),
4 TOP(20,20),BOT(20),UP(20,20),DOWN(20)
DIMENSION TITLE(20),IFMT(20),VAR(20)

INPUT SECTION: READ IN TITLE CARD, PROBLEM CARD, OPTIONAL CARDS,
FORMAT CARD, AND DATA CARDS
WRITE OUT SOLUTION SPECIFICATIONS

IFINE = 0
READ (5,900) (TITLE(I),I=1,20)
WRITE (6,903) (TITLE(I),I=1,20)
READ (5,901) NOBS,NVARS,NGPS,ICRIT,NOSTAN,IDIST,NODATA,KTIME
IF (NOBS.EQ.0) GO TO 820
IF (NOBS.GT.600) GO TO 815
IF (NVARS.GT.20) GO TO 815
IF (NGPS.GT.20) GO TO 815
IF (ICRIT.GT.4) GO TO 815
IF (NOSTAN.GT.1) GO TO 815
IF (IDIST.GT.2) GO TO 815
IF (NVARS.EQ.1) ICRIT=1
IF (NVARS.EQ.1) IDIST=0
IF (NODATA.EQ.1) GO TO 5
4 READ (5,902) (IFMT(I),I=1,18)
5 WRITE (6,904) NOBS,NVARS,NGPS
GO TO 10,20,30,40,ICRIT
10 WRITE (6,920)
GO TO 50
20 WRITE (6,921)
GO TO 50
30 WRITE (6,922)
GO TO 50
40 WRITE (6,923)
50 CONTINUE
60 IF (NOSTAN) 60,60,70
WRITE (6,924)
GO TO 80
70 WRITE (6,925)
80 CONTINUE
85 IF (IDIST) 85,85,90
WRITE (6,926)
GO TO 95
90 IF (IDIST.EQ.2) GO TO 92
WRITE (6,936)
GO TO 95
92 WRITE (6,927)

```

```

RRR01770
RRR01780
RRR01790
RRR01800
RRR01810
RRR01820
RRR01830
RRR01840
RRR01850
RRR01860
RRR01870
RRR01880
RRR01890
RRR01900
RRR01910
RRR01920
RRR01930
RRR01940
RRR01950
RRR01960
RRR01970
RRR01980
RRR01990
RRR02000
RRR02010
RRR02020
RRR02030
RRR02040
RRR02050
RRR02060
RRR02070
RRR02080
RRR02090
RRR02100
RRR02110
RRR02120
RRR02130
RRR02140
RRR02150
RRR02160
RRR02170
RRR02180
RRR02190
RRR02200
RRR02210
RRR02220
RRR02230

```

CC

```

55 CONTINUE
   IF (NODATA.EQ.1) GO TO 101
   WRITE (6,905) (IFMT(I),I=1,13)
   DO 100 I = 1, NODATA
   READ (5,IFMT) IDENT(I), (DATA(I,J),J=1,NVARS)
100 CONTINUE = 1
   GO TO 102
101 WRITE (6,937)
102 WRITE (6,943) KTIME
C
C
C
      CALCULATE VARIABLE MEANS AND VARIANCES
      SUBTRACT OUT OVERALL MEAN
      DO 105 J=1, NVARS
      VMEAN (J) = 0.0
      VAR (J) = 0.0
      DO 110 I = 1, NODATA
      DC 110 J = 1, NVARS
      VMEAN(J)=VMEAN(J)+NS(I)*DATA(I,J)/NODATA
      DO 120 I=1, NODATA
      DO 120 J=1, NVARS
      DATA(I,J)=DATA(I,J)-VMEAN(J)
      VAR(J)=VAR(J)+NS(I)*(DATA(I,J)**2)/(NODATA-1)
120
      DO 125 J = 1, NVARS
      IF (VAR(J).LE.0.000001) MFLAG=1
125 SD(J) = SQRT (VAR(J))
      IF (MFLAG.EQ.1) GO TO 825
C
C
      CALCULATE T = X'X, THE CROSS-PRODUCTS MATRIX
131 DO 135 K=1, NVARS
135 T(K,J) = 0.0
      DO 140 K = 1, NVARS
      DO 140 J=K, NVARS
      T(K,J)=T(K,J)+NS(I)*DATA(I,K)*DATA(I,J)
140
      OUTPUT SECTION: WRITE OUT MEANS, STANDARD DEVIATIONS, AND CROSS-
      PRODUCTS MATRIX T
      WRITE (6,908)
      WRITE (6,907) (VMEAN(J),J=1,NVARS)
      WRITE (6,909)
      WRITE (6,907) (SD(J),J=1,NVARS)
      WRITE (6,911)
      DC 810 I = 1, NVARS

```

```

RRR02240
RRR02250
RRR02260
RRR02270
RRR02280
RRR02290
RRR02300
RRR02310
RRR02320
RRR02330
RRR02340
RRR02350
RRR02360
RRR02370
RRR02380
RRR02390
RRR02400
RRR02410
RRR02420

```

```

RRR02450
RRR02460
RRR02480
RRR02490
RRR02500
RRR02510
RRR02520
RRR02530
RRR02540
RRR02550
RRR02560
RRR02570
RRR02580
RRR02590
RRR02600

```

```

RRR02620
RRR02630
RRR02640
RRR02650
RRR02660
RRR02670
RRR02680
RRR02690
RRR02700
RRR02710

```

```

      WRITE (6,907) (T(J,I),J=1,I)
      810 CONTINUE
      IF (NOSTAN.EQ.0) GO TO 180
C
C      STANDARDIZE IF REQUESTED
C
      DO 160 J=1,NVARS
      DO 150 I=1,NOBS
      DATA(I,J)=DATA(I,J) / SD(J)
      150 DO 160 K=J,NVARS
      T(J,K) = (1.0/SD(J))*T(J,K)*(1.0/SD(K))
      160 RETURN
      815 IFINE = 1
      IF (NODATA.EQ.0) IFINE=2
      WRITE (6,935)
      RETURN
      820 IFINE = 2
      RETURN
      825 WRITE (6,945) J
      IFINE = 1
      RETURN
C
C      FCRMAT STATEMENTS
C
      900 FORMAT (20A4)
      901 FORMAT (14,2I2,5I1)
      902 FORMAT (18A4)
      903 FORMAT (1,1,20A4)
      904 FCRMAT (0, THE NUMBER OF OBSERVATIONS IS ,14,/, THE NUMBER OF
      1 VARIABLES IS ,12,/, THE NUMBER OF GROUPS (INITIAL) IS ,12)
      905 FCRMAT (0, THE INPUT FORMAT IS ,18A4)
      906 FCRMAT (0,10(F11.3,1X))
      907 FCRMAT (0, THE VARIABLE MEANS ARE )
      908 FCRMAT (0, THE VARIABLE STANDARD DEVIATIONS ARE )
      909 FCRMAT (0, THE CROSS-PRODUCTS MATRIX IS )
      910 FCRMAT (0, THE CRITERION IS TO MINIMIZE THE TRACE W )
      911 FCRMAT (0, THE CRITERION IS TO MINIMIZE THE DETERMINANT W )
      920 FCRMAT (0, THE CRITERION IS TO MAXIMIZE THE LARGEST ROOT OF W-1B )
      921 FCRMAT (0, THE CRITERION IS TO MAXIMIZE THE TRACE OF W-1B )
      922 FCRMAT (0, THE CRITERION IS TO MINIMIZE THE TRACE OF W-1B )
      923 FCRMAT (0, THE CRITERION IS TO MINIMIZE THE TRACE OF W-1B )
      924 FCRMAT (0, THE CRITERION IS TO MINIMIZE THE TRACE OF W-1B )
      925 FCRMAT (0, THE CRITERION IS TO MINIMIZE THE TRACE OF W-1B )
      926 FCRMAT (0, THE CRITERION IS TO MINIMIZE THE TRACE OF W-1B )
      927 FCRMAT (0, THE CRITERION IS TO MINIMIZE THE TRACE OF W-1B )
      928 FCRMAT (0, THE CRITERION IS TO MINIMIZE THE TRACE OF W-1B )
      929 FCRMAT (0, THE CRITERION IS TO MINIMIZE THE TRACE OF W-1B )
      930 FCRMAT (0, THE CRITERION IS TO MINIMIZE THE TRACE OF W-1B )
      931 FCRMAT (0, THE CRITERION IS TO MINIMIZE THE TRACE OF W-1B )
      932 FCRMAT (0, THE CRITERION IS TO MINIMIZE THE TRACE OF W-1B )
      933 FCRMAT (0, THE CRITERION IS TO MINIMIZE THE TRACE OF W-1B )
      934 FCRMAT (0, THE CRITERION IS TO MINIMIZE THE TRACE OF W-1B )
      935 FCRMAT (0, THE CRITERION IS TO MINIMIZE THE TRACE OF W-1B )
      936 FCRMAT (0, THE CRITERION IS TO MINIMIZE THE TRACE OF W-1B )
      937 FCRMAT (0, THE CRITERION IS TO MINIMIZE THE TRACE OF W-1B )
      943 FCRMAT (0, THE CRITERION IS TO MINIMIZE THE TRACE OF W-1B )
      945 FCRMAT (0, THE CRITERION IS TO MINIMIZE THE TRACE OF W-1B )

```



```

210 YVEC(J) = SVCENT(M,J)
    CALL DISTCE (XVEC,YVEC,DISTA)
    IF (IFINE.EQ.1) RETURN
    IF (M.EQ.1) GO TO 212
    IF (DISTA.GE.SMDIST) GO TO 215
212 SMDIST = DISTA
    IDATAT(I) = M
215 ASSIGN OBSERVATION TO CLOSEST CLUSTER AND UPDATE THAT CLUSTER
    CENTER
    K = IDATAT(I)
    NISVT(K) = NISVT(K) + 1
    BOT(K) = BOT(K) + NS(I)
    DO 225 J = 1, NVARS
    TOP(K,J) = TOP(K,J) + NS(I)*DATA(I,J)
    SVCENT(K,J) = TOP(K,J)/BOT(K)
    IDIST = ITEMP
    END OF DO FOR EACH OBSERVATION
    RECALCULATE CLUSTER CENTERS TO ELIMINATE INITIAL RANDOM OBSERVA-
    TIONS
225 DO 230 M=1,NGPS
    BOT(M) = 0.0
    NISVT(M) = 1
    DO 230 J=1, NVARS
    TOP(M,J) = 0.0
    DC 235 I=1, NOBS
    M = IDATAT(I)
    BOT(M) = BOT(M) + NS(I)
    DO 235 J=1, NVARS
    TOP(M,J) = TOP(M,J) + NS(I)*DATA(I,J)
    SVCENT(M,J) = TOP(M,J)/BOT(M)
230
235
    CALCULATE THE B AND W MATRICES
    CALCULATE THE CRITERION VALUE
    IDIR=2
    IF (ICRIT.EQ.1) IDIR=1
    CALL WCALC (SVCENT,NISVT,NGPST,IDIR)
    IF (IFINE.EQ.1) RETURN
    CALL CRITON (CRIT)
    IF (IFINE.EQ.1) RETURN
    WHICH INITIAL CONFIGURATION
    IF (IBSRT.GT.1) GO TO 250
    FIRST INITIAL CONFIGURATION: BCRIT IS THE BEST CRITERION
    245 BCRIT = CRIT
    SAVE CLUSTER SIZES (NISV), CLUSTER CENTERS (SVCENT), CLUSTER
    LISTS (LSTSV), AND OBSERVATION ID'S (IDATA)

```

```

RRR03620
RRR03630
RRR03640
RRR03650
RRR03660
RRR03670
RRR03680
RRR03690
RRR03700
RRR03710
RRR03720
RRR03730
RRR03740

RRR03760
RRR03770
RRR03780
RRR03790
RRR03800

RRR03810
RRR03820

RRR03830
RRR03840
RRR03850

RRR03860

RRR03880
RRR03890
RRR03900
RRR03910
RRR03920
RRR03930
RRR03940
RRR03950
RRR03960
RRR03970
RRR03980
RRR03990
RRR04000
RRR04010
RRR04020
RRR04030

```

```

DO 247 M = 1,NGPS
NISV(M) = NISVT(M)
DOWN(M) = BOT(M)
DO 247 L = 1,NVARS
UP(M,L) = TOP(M,L)
247 SVCEN(M,L) = SVCENT(M,L)
DO 249 I = 1,NOBS
249 ILATA(I) = IDATAT(I)
GO TO 260
C SECOND OR THIRD INITIAL CONFIGURATION
250 IF (CRIT.GE.BCRIT) GO TO 260
255 BCRIT = CRIT
DO 257 M = 1,NGPS
NISV(M) = NISVT(M)
DOWN(M) = BOT(M)
DO 257 L = 1,NVARS
UP(M,L) = TOP(M,L)
257 SVCEN(M,L) = SVCENT(M,L)
DO 259 I = 1,NOBS
259 IDATA(I) = IDATAT(I)
260 CONTINUE
C ALL DONE
C RETURN
C END
SUBROUTINE KMEANS (CRIT)
C K-MEANS. EACH OBSERVATION IS ASSIGNED TO THE CLOSEST CLUSTER
C CENTER.
COMMON NOBS,NVARS,NGPS,ICRIT,NOSTAN,IDIST,IFIME,KTIME,IDENT(600),
1 IDATA(600,20),NISV(20),VMEAN(20),SD(20),XVEC(20),YVEC(20),BT(20,20),
2 NISVT(20),SVCENT(20,20),IDATAT(600),VEC(20,20),EIG(20),NS(600),
3 TOP(20,20),BOT(20),UP(20,20),DOWN(20)
C INITIALIZING FLAGS
C AND TEMPORARY STORAGE FOR NISV, SVCEN
200 NGPST = NGPS
JFLAG = 0
BCRIT = CRIT
DO 201 M = 1,NGPS
NISVT(M) = NISV(M)
BOT(M) = DOWN(M)
DO 201 J = 1,NVARS
TOP(M,J) = UP(M,J)
201 SVCENT(M,J) = SVCEN(M,J)
IDIR = 1
C RECALCULATING WITHIN-CLUSTERS MATRIX
RRR04040
RRR04050
RRR04060
RRR04070
RRR04080
RRR04090
RRR04100
RRR04110
RRR04120
RRR04130
RRR04140
RRR04150
RRR04160
RRR04170
RRR04180
RRR04190
RRR04200
RRR04210
RRR04220
RRR04230
RRR04240
RRR04250
RRR04260
RRR04270
RRR04280
RRR04290
RRR04300
RRR04310
RRR04320
RRR04340
RRR04350
RRR04360
RRR04370
RRR04380
RRR04390
RRR04400
RRR04410
RRR04420
RRR04430
RRR04440

```

```

IF (IDIST.EQ.2) IDIR=2
CALL WCALC (SVCENT,NISVT,NGPST,IDIR)
IF (IFINE.EQ.1) RETURN
INITIALIZING TEMPORARY STORAGE FOR LSTSV, IDATA
DO 210 I = 1,NBRS
210 IDATAT(I) = IDATA(I)
C MAJOR DO LOOP: DO FOR EACH OBSERVATION
DO 400 I = 1,NBRS
C
C BEGIN K-MEANS SECTION
C
C CALCULATE VECTOR OF DISTANCES FROM OBSERVATION TO EACH CLUSTER
CENTER
325 DC 335 M=1,NGPST
DO 330 J=1,NVARS
XVEC(J) = DATA(I,J)
330 YVEC(J) = SVCENT(M,J)
CALL DISTCE(XVEC,YVEC,DISTB)
IF (IFINE.EQ.1) RETURN
IF (M.EQ.1) GO TO 332
IF (DISTB.GE.SMDIST) GO TO 335
332 SMDIST = DISTB
335 IDGP = M
335 CONTINUE
C IS OBSERVATION ALREADY IN THE CLOSEST CLUSTER? IF YES, SKIP THIS
SECTION
IF (IDGP.EQ.IDATAT(I)) GO TO 360
346 IOLD IS OLD CLUSTER ASSIGNMENT
IOLD = IDATAT(I)
C INEW IS NEW CLUSTER ASSIGNMENT
INEW = IDGP
JFLAG = 1
C RECALCULATE CLUSTER CENTERS
BOT(IOLD) = BOT(IOLD) - NS(I)
BOT(INEW) = BOT(INEW) + NS(I)
DC 348 J = 1,NVARS
TOP(IOLD,J) = TOP(IOLD,J) - NS(I)*DATA(I,J)
SVCENT(IOLD,J) = TOP(IOLD,J)/BOT(IOLD)
TOP(INEW,J) = TOP(INEW,J) + NS(I)*DATA(I,J)
SVCENT(INEW,J) = TOP(INEW,J)/BOT(INEW)
348 ADJUST NISV
NISVT(IOLD) = NISVT(IOLD)-1
NISVT(INEW) = NISVT(INEW)+1
ADJUST IDATA
IDATAT(I) = IDGP
C RECALCULATE WITHIN-CLUSTERS MATRIX FOR USE IN COMPUTING SCALED
EUCLIDIAN AND MAHALANOBIS DISTANCE
IF (IDIST.EQ.0) GO TO 360

```

RRR04450
RRR04460
RRR04470
RRR04480
RRR04490
RRR04500
RRR04510
RRR04520
RRR04530
RRR04540
RRR04550
RRR04560
RRR04570
RRR04580
RRR04590
RRR04600
RRR04610
RRR04620
RRR04630
RRR04640
RRR04650
RRR04660
RRR04670
RRR04680
RRR04690
RRR04700
RRR04710
RRR04720
RRR04730
RRR04740
RRR04750
RRR04760
RRR04770

RRR04780

RRR04830
RRR04840
RRR04850
RRR04860
RRR04870
RRR04880
RRR04890
RRR04900

```

CALL WCALC (SVCENT,NISVT,NGPST,IDIR)
IF (IFINE.EC.1) RETURN
360 CONTINUE
C   DONE WITH MAJOR DO LOOP
400 CONTINUE
C
C   CALCULATE THE CRITERION
C
RECALCULATE SVCEN:  ACCURACY MEASURE
DC 401 M=1,NGPST
BOT(M)=0.0
DO 401 J=1,NVARS
TOP(M,J)=0.0
SVCENT(M,J)=0.0
401 DC 402 I=1,NOBS
M = IDATAT(I)
BOT(M)=BOT(M)+NS(I)
DC 402 J=1,NVARS
TOP(M,J)=TOP(M,J)+NS(I)*DATA(I,J)
SVCENT(M,J)=TOP(M,J)/BOT(M)
402 RECALCULATE WITHIN-CLUSTERS MATRIX AND CRITERION VALUE BASED ON
C   NEW CLUSTER CENTER VALUES
C   IDIR = 2
IF (ICRIT.EQ.1) IDIR = 1
CALL WCALC (SVCENT,NISVT,NGPST,IDIR)
IF (IFINE.EC.1) RETURN
405 CALL CRITON (CRIT)
IF (IFINE.EQ.1) RETURN
C   IF CRITERION BETTER THAN BEFORE? IF YES, THEN ANOTHER ITERATION:
C   IF NO, FINISH KMEANS
C   IF (CRIT.GE.BCRIT) GO TO 535
C
C   ANOTHER ITERATION
C
PUT TEMPORARY VALUES INTO PERMANENT LOCATIONS
515 NGPS=NGPST
DC 520 M=1,NGPS
DOWN(M)=BOT(M)
NISV(M)=NISVT(M)
DO 520 J=1,NVARS
UP(M,J)=TOP(M,J)
520 SVCEN(M,J)=SVCENT(M,J)
DO 530 I=1,NOBS
530 IDATA(I)=IDATAT(I)
GO TO 200
C
C   FINISH
C

```

```

RRR04910
RRR04920
RRR04930
RRR04940
RRR04950
RRR04960
RRR04970
RRR04980
RRR04990
RRR05000
RRR05010
RRR05020
RRR05030
RRR05040
RRR05050
RRR05070
RRR05080
RRR05090
RRR05100
RRR05110
RRR05120
RRR05130
RRR05140
RRR05150
RRR05160
RRR05170
RRR05180
RRR05190
RRR05200
RRR05210
RRR05220
RRR05230
RRR05240
RRR05250
RRR05260
RRR05270
RRR05280
RRR05290
RRR05300
RRR05310
RRR05320

```



```

C      JFLAG = 0 MEANS NO CHANGES HAVE BEEN MADE DURING THE LAST
C      ITERATION: ITERATIONS CONVERGED
C 535 IF (JFLAG.EQ.0) RETURN
C      JFLAG = 1 MEANS CHANGES HAVE BEEN MADE BUT THE CRITERION VALUE
C      GOT WORSE: ITERATIONS NOT CONVERGING
C 540 WRITE (6,940)
C      WRITE (6,942) CRIT
C      WRITE (6,943) BCRTT
C      RECALCULATE WITHIN-CLUSTERS MATRIX AND RESET CRIT
C      IDIR=2
C      IF (ICRIT.EQ.1) IDIR=1
C      CALL WCALC (SVCEN,NISV,NGPS,IDIR)
C      IF (IFINE.EQ.1) RETURN
C      CRIT=BCRTT
C      RETURN
C 940 FORMAT ('0 ITERATIONS NOT CONVERGING')
C 942 FORMAT ('0 THE CRITERION VALUE IS ',E12.6)
C 943 FORMAT ('0 THE BEST CRITERION VALUE IS ',E12.6)
C      END
C      SUBROUTINE ISWITCH (CRIT)
C
C      THIS SUBROUTINE CONSIDERS SWITCHING EACH OBSERVATION TO A
C      DIFFERENT CLUSTER. THE SWITCH IS MADE IFF A BETTER CRITERION
C      VALUE RESULTS.
C      THIS SUBROUTINE ALSO DEPENDS ON THE PARAMETER KTIME: IT DETERMINE
C      WHICH OBSERVATIONS ARE TO BE CONSIDERED. FOR COMPLETE EXPLANATION
C      OF THE KTIME PARAMETER, SEE THE PROGRAM DESCRIPTION.
C
C      COMMON NOBS,NVARS,NGPS,ICRIT,NOSTAN,IDIST,IFINE,KTIME,IDENT(600),
C      1DATA(600),NISV(20),W(20,20),WFC(20,20),SVCEN(20,20),
C      21DATA(600),NISV(20),VMEAN(20),SD(20),XVEC(20),YVEC(20),RT(20,20),
C      3NISVT(20),SVCENT(20,20),IDATA(600),VEC(20,20),EIG(20),NS(600),
C      4TOP(20,20),BOT(20),UP(20,20),DOWN(20)
C      DIMENSION JFLAG (20)
C      SET UP THE IDIR FLAG
C      IDIR=2
C      IF (ICRIT.EQ.1) IDIR=1
C      KTIME = 0 MEANS SKIP THIS HEURISTIC
C      IF (KTIME.EC.0) GO TO 750
C      SET UP TEMPORARY STORAGE AREAS FOR NGPS, NISV, AND SVCEN
C      NGPST=NGPS
C      DO 500 M = 1,NGPS
C      BOT(M)=DOWN(M)
C      NISVT(M)=NISV(M)
C      JFLAG(M)=0
C      DO 500 J = 1,NVARS
C      TOP(M,J)=UP(M,J)
C      SVCENT(M,J)=SVCEN(M,J)
C 500

```

```

RRR055330
RRR055340
RRR055350
RRR055360
RRR055370
RRR055380
RRR055390
RRR055400
RRR055410
RRR055420
RRR055430
RRR055440
RRR055450
RRR055460
RRR055470
RRR055480
RRR055490
RRR055500
RRR055510
RRR055520
RRR055530
RRR055540
RRR055550
RRR055560
RRR055570
RRR055580
RRR055590
RRR055600
RRR055610
RRR055620
RRR055630

RRR05650
RRR05660
RRR05670
RRR05680
RRR05690
RRR05700
RRR05710
RRR05720
RRR05730

RRR05740
RRR05750
RRR05760
RRR05770

```

```

C      SET IFLAG = 0:  WILL RE SET = 1 IF ANY SWITCH IS MADE
60C IFLAG=0
C      MAJOR DO LOOP:  DO 700
C      DO 700 M=1, NGPS
C      JFLAG(M).EQ.2) GO TO 700
C      IF (JFLAG(M).EQ.2) GO TO 700
C      KTIME = 9 MEANS ALL OBSERVATIONS WILL BE CONSIDERED
C      IF (KTIME.EQ.9) GO TO 617
C      COMPUTE ALL DISTANCES FROM OBSERVATIONS TO CLUSTER CENTERS FOR THE
C      OBSERVATIONS IN CLUSTER M
C      DO 605 J=1, NVARS
C      605 YVEC(J) = $VCEN(M,J)
C      BIGDIS = 0.
C      DO 615 I=1, NOBS
C      IF (I.DAT(I)).NE.M) GO TO 615
C      DO 610 J=1, NVARS
C      610 XVEC(J) = DATA(I,J)
C      CALL DISTCE (XVEC,YVEC,DISTA)
C      LOCATE OBSERVATION WITH BIGGEST DISTANCE
C      IF (BIGDIS.GE.DISTA) GO TO 615
C      BIGDIS = DISTA
C      615 CONTINUE
C      DIVIDE BIGDIS BY TIMING PARAMETER KTIME
C      TIME = KTIME
C      PSME = BIGDIS / TIME
C      617 MOLD = M
C      DO FOR ALL OBSERVATIONS IN CLUSTER M
C      DO 691 I=1, NOBS
C      IF (I.DAT(I)).NE.M) GO TO 691
C      KFLAG = 0
C      KTIME = 9 MEANS CONSIDER ALL OBSERVATIONS
C      IF (KTIME.EQ.9) GO TO 618
C      DO NO CONSIDER SWITCHING THE OBSERVATION IF THE DISTANCE TO ITS
C      CURRENT CLUSTER CENTER IS LT BIGDIS / KTIME
C      DO 619 J=1, NVARS
C      619 XVEC(J) = DATA(I,J)
C      CALL DISTCE (XVEC,YVEC,DISTA)
C      IF (DISTA.LT.PSME) GO TO 691
C      618 CONTINUE
C      DO 690 MNEW=1, NGPS
C      IF (NISV(MOLD).EQ.1) GO TO 690
C      IF (MNEW.EQ.MOLD) GO TO 690
C      IF (KFLAG.EQ.1) GO TO 690
C      COMPUTE NEW SVCEM BASED ON OBSERVATION BEING SWITCHED
C      BOT(MNEW)=BOT(MNEW)+NS(I)
C      BOT(MOLD)=BOT(MOLD)-NS(I)
C      DO 620 J=1, NVARS
C      TOP(MNEW,J)=TOP(MNEW,J)+NS(I)*DATA(I,J)
C      620 J=1, NVARS
C      6220

```

```

RRR05780
RRR05790
RRR05800
RRR05810
RRR05820
RRR05830
RRR05840
RRR05850
RRR05860
RRR05870
RRR05880
RRR05890
RRR05900
RRR05910
RRR05920
RRR05930
RRR05940
RRR05950
RRR05960
RRR05970
RRR05980
RRR05990
RRR06000
RRR06010
RRR06020
RRR06030
RRR06040
RRR06050
RRR06060
RRR06070
RRR06080
RRR06090
RRR06100
RRR06110
RRR06120
RRR06130
RRR06140
RRR06150
RRR06160
RRR06170
RRR06180
RRR06190
RRR06200
RRR06210
RRR06220

```

```

        SVCENT(MNEW,J)=TOP(MNEW,J)/BOT(MNEW)
        TCP(MOLD,J)=TOP(MOLD,J)-NS(I)*DATA(I,J)
        SVCENT(MOLD,J)=TOP(MOLD,J)/BOT(MOLD)
620  CONTINUE
      C  ADJUST NISV
        NISVT(MNEW) = NISV(MNEW)+1
        NISVT(MOLD) = NISV(MOLD)-1
      C  COMPUTE WITHIN-CLUSTERS MATRIX AND NEW CRITERION VALUE
        CALL WCALC (SVCENT,NISVT,NGPST,DIR)
        IF (IFINE.EC.1) RETURN
        CALL CRITON (TCRIT)
        IF (IFINE.EQ.1) RETURN
      C  MAKE THE SWITCH IFF NEW CRIT IS BETTER
        MAKE (TCRIT.GE.CRIT) GO TO 680
      C  MAKE THE SWITCH
      C
      C  RESET FLAGS
        JFLAG(M) = 1
        JFLAG(MNEW) = 1
        IFLAG=1
      C  KFLAG = 1
        ADJUST CRIT, IDATA, NISV, SVCEN
        CRIT=TCRIT
        ICATA(I) = MNEW
        DCWN(MNEW) = BOT(MNEW)
        DCWN(MOLD) = BOT(MOLD)
        NISV(MNEW) = NISVT(MNEW)
        NISV(MOLD) = NISVT(MOLD)
629  DO 630 J=1,NVARS
        UP(MNEW,J)=TOP(MNEW,J)
        UP(MOLD,J)=TCP(MOLD,J)
        SVCEN(MNEW,J) = SVCENT(MNEW,J)
        SVCEN(MOLD,J) = SVCENT(MOLD,J)
630  GO TO 690
      C  SWITCH WAS NOT MADE; RESET NISVT, SVCENT, TO VALUES PRESENT
        BEFORE THE SWITCH WAS CONSIDERED
      C 680 NISVT(MNEW) = NISV(MNEW)
        NISVT(MOLD) = NISV(MOLD)
        BOT(MNEW)=DCWN(MNEW)
        BOT(MOLD)=DCWN(MOLD)
        DO 685 J=1,NVARS
        TOP(MNEW,J)=UP(MNEW,J)
        TOP(MOLD,J)=UP(MOLD,J)
        SVCENT(MNEW,J) = SVCEN(MNEW,J)
        SVCENT(MOLD,J) = SVCEN(MOLD,J)
685  CONTINUE
690  CONTINUE
691

```

PRR06270
 PRR06280
 PRR06290
 PRR06300
 PRR06310
 PRR06320
 PRR06330
 PRR06340
 PRR06350
 PRR06360
 PRR06370
 PRR06380
 PRR06390
 PRR06400
 PRR06410
 PRR06420
 PRR06430
 PRR06440
 PRR06450
 PRR06460
 PRR06470
 PRR06480

PRR06490
 PRR06500
 PRR06510

PRR06520
 PRR06530
 PRR06540
 PRR06550
 PRR06560
 PRR06570
 PRR06580

PRR06590

PRR06600
 PRR06610
 PRR06620
 PRR06630

```

C      FINISH WITH CLUSTER M: IF NO SWITCHES HAVE BEEN MADE, SET
C      JFLAG(M)=2 AND GO TO NEXT CLUSTER: IF SWITCHES HAVE BEEN MADE,
C      ADJUST LSTSV AND SET JFLAG(M)=0 AND GO TO NEXT CLUSTER
695 IF (JFLAG(M).EQ.0) GO TO 699
      JFLAG(M)=0
      GO TO 700
699 JFLAG(M)=2
700 CONTINUE
C      DONE WITH ALL CLUSTERS: IFLAG=1 MEANS SOME SWITCHES HAVE BEEN
C      MADE: GO BACK AND ITERATE
C      IF (IFLAG.EQ.1) GO TO 600
C      ALL DONE: ACCURATELY CALCULATE MEANS AND CRITERION AND OUTPUT
C      THE RESULTS
750 WRITE (6,900)
900 FORMAT (10 THE FINAL CLUSTER SOLUTION IS ')
C      RECALCULATE CLUSTER CENTERS, WITHIN-CLUSTERS MATRIX, AND CRITERION
C      VALUE
      DO 715 M=1,NGPS
        BCT(M)=0.0
      DO 715 J=1,NVARS
        TCPP(M,J)=0.0
      SVCEN(M,J)=0.0
      DO 720 I=1,NOBS
        M=IDATA(I)
        BCT(M)=BCT(M)+NS(I)
      DO 720 J=1,NVARS
        TOP(M,J)=TOP(M,J)+NS(I)*DATA(I,J)
      SVCEN(M,J)=TOP(M,J)/BCT(M)
      CALL WCALC (SVCEN,NISV,NGPS,IDIR)
      IF (IFINE.EC.1) RETURN
      CALL CRITON (CRIT)
      IF (IFINE.EQ.1) RETURN
      CALL THE OUTPUT ROUTINE
      CALL OUTPUT (CRIT)
      WRITE (6,901)
901 FORMAT (10 END OF CLUSTER PROBLEM')
      RETURN
      END
      SUBROUTINE OUTPUT (CRIT)
C      THIS SUBROUTINE PRINTS OUT THE CLUSTER SOLUTION
C      FOR EACH CLUSTER - THE CLUSTER SIZE
C      THE CLUSTER CENTROID
C      THE OBSERVATIONS BELONGING TO THE CLUSTER
C      THE WITHIN CELLS MATRIX
C      THE CRITERION VALUE

```

```

RRR06640
RRR06650
RRR06660
RRR06670
RRR06680
RRR06690
RRR06700
RRR06710
RRR06720
RRR06730
RRR06740
RRR06750
RRR06760
RRR06770
RRR06780
RRR06790
RRR06800
RRR06810
RRR06820
RRR06830

RRR06840

RRR06850
RRR06860
RRR06870

RRR06880

RRR06900
RRR06910
RRR06920
RRR06930
RRR06940
RRR06950
RRR06960
RRR06970
RRR06980
RRR06990
RRR07000
RRR07010
RRR07020
RRR07030
RRR07040
RRR07050
RRR07060
RRR07070

```



```

C
COMMON NOBS,NVARS,NGPS,ICRIT,NOSTAN,IDIST,IFINE,KTIME,IDENT(6CC),
1DATA(600),T(20,20),B(20,20),W(20,20),WFCT(20,20),SVCEN(20,20),
2IDATA(600),NISV(20),VMEAN(20),SD(20),XVEC(20),YVEC(20),BT(20,20),
3NISVT(20),SVCENT(20,20),IDATAT(600),VEC(20,20),EIG(20),NS(600),
4TOP(20,20),BOT(20),UP(20,20),DOWN(20)
UNSTANDARDIZE IF NECESSARY
IF (NOSTAN.EQ.0) GO TO 40
DO 20 M=1,NGPS
DO 20 J=1,NVARS
20 SVCEN(M,J) = SVCEN(M,J)*SD(J)
DO 30 K=1,NVARS
30 T(J,K) = SD(J)*T(J,K)*SD(K)
IF IDIR = 2
IF (ICRIT.EQ.1) IDIR=1
CALL WCALC (SVCEN,NISV,NGPS,IDIR)
IF (IFINE.EQ.1) RETURN
CALL CRITON (CRIT)
IF (IFINE.EQ.1) RETURN
ADD THE OVERALL MEAN BACK INTO EACH OBSERVATION: UNSTANDARDIZE
C
C
40 DO 80 J=1,NVARS
DO 50 I=1,NCBS
IF (NOSTAN.EQ.1) DATA(I,J) = DATA(I,J)*SD(J)
50 DATA(I,J) = DATA(I,J) + VMEAN(J)
DO 60 M=1,NGPS
60 SVCEN(M,J) = SVCEN(M,J) + VMEAN(J)
80 CONTINUE
C
WRITE OUT NISV, SVCEN
DO 115 M=1,NGPS
WRITE (6,900) M
WRITE (6,901) NISV(M)
WRITE (6,902) (SVCEN(M,J),J=1,NVARS)
WRITE (7,915) (SVCEN(M,J),J=1,NVARS)
FORMAT(8F8.4,/8F8.4)
WRITE (6,903)
LSTNO = NISV(M)
K = 0
C
DO 100 I=1,NOBS
IF (IDATA(I).NE.M) GO TO 100
K = K+1
IDATAT(K) = IDENT(I)
100 CONTINUE
C
SCRT THE OBSERVATION IDENTIFICATIONS
L=0
101 L=L+1
IF (L.EQ.NISV(M)) GO TO 110

```

RRR07080
RRR07090
RRR07100
RRR07110

RRR07130
RRR07140
RRR07150
RRR07160
RRR07170
RRR07180
RRR07190
RRR07200
RRR07210
RRR07220
RRR07230
RRR07240
RRR07250
RRR07260
RRR07270
RRR07280
RRR07290
RRR07300
RRR07310
RRR07320
RRR07330
RRR07340
RRR07350
RRR07360
RRR07370
RRR07380
RRR07390
RRR07400

RRR07410
RRR07420
RRR07430
RRR07440
RRR07450
RRR07460
RRR07470
RRR07480
RRR07490
RRR07500
RRR07510
RRR07520

```

RRR07530
RRR07540
RRR07550
RRR07560
RRR07570
RRR07580
RRR07590
RRR07600
RRR07610
RRR07620
RRR07630

```

```

RRR07640
RRR07650
RRR07660
RRR07670
RRR07680
RRR07690
RRR07700
RRR07710
RRR07720
RRR07730
RRR07740
RRR07750
RRR07760
RRR07770
RRR07780
RRR07790
RRR07800
RRR07810
RRR07820
RRR07830
RRR07840
RRR07850
RRR07860
RRR07870
RRR07880
RRR07890
RRR07900
RRR07910
RRR07920
RRR07930
RRR07940
RRR07950
RRR07960
RRR07970
RRR07980

```

```

102 LL=L+1
   IF (IDATAT(L).LE.IDATAT(LL)) GO TO 101
   LTEMP = IDATAT(L)
   IDATAT(L) = IDATAT(LL)
   IDATAT(LL) = LTEMP
   IF (L.EQ.1) GO TO 102
   L=L-1
   GO TO 102
110 CONTINUE
   WRITE OUT THE IDENTIFICATIONS
   WRITE (6,904) (IDATAT(K),K=1,LSINC)
   WRITE (7,920) (IDATAT(K),K=1,LSINC)
   FCRMAT(1615)
115 CONTINUE
   WRITE OUT THE WITHIN-CLUSTERS MATRIX
119 WRITE (6,905)
   DO 120 J=1,NVARS
     WRITE (6,906) (W(K,J),K=1,J)
120 CONTINUE
   ICRIT IS THE CRITERION CHOSEN BY THE USER
   GO TO (150,160,170,178),ICRIT
   WRITE OUT TRACE W
150 WRITE (6,907) CRIT
   RETURN
   WRITE OUT DET W
   ANL LOG (DET T / DET W)
160 CALL UPFCT (NVARS,T,M)
   DET = 1.
165 DC 165 J=1,NVARS
   DET = DET*(J,J)
   WRITE (6,912) CRIT
   CRIT = ALOG10 ((DET**2)/(CRIT**2))
   WRITE (6,909) CRIT
   RETURN
   FOR LARGEST ROOT AND HOTELLING'S TRACE CRITERIA, FIND EIGENVALUES
   FOR IB-GW=0.
170 DO 175 J=1,NVARS
   DO 175 K=J,NVARS
     BT(J,K) = B(J,K)
     BT(K,J) = BT(J,K)
     CRIT = 1.0 / CRIT
     CALL UTISU (NVARS,BT,WFCF)
     CALL EIGN (NVARS,BT,EIG,VEC,IND)
     IF (IND.NE.0) GO TO 180
175 WRITE OUT EIGENVALUES
   WRITE (6,910) (EIG(J),J=1,NVARS)
   IF (ICRIT.EQ.4) GO TO 178
   WRITE OUT LARGEST ROOT

```


134

```

C      DC 100 J=1,NVARS
      DO 100 K=J,NVARS
      DO 105 M=1,NGP
      DO 105 J=1,NVARS
      DO 105 K=J,NVARS
      BT(J,K) = SV(M,J) * SV(M,K)
      E(J,K) = B(J,K)+BT(J,K)*BOT(M)
      DO 110 J=1,NVARS
      DO 110 K=J,NVARS
      W(J,K) = T(J,K) - B(J,K)
      C      CALCULATE CHOLESKY FACTOR OF W IF IDIR = 2
      C      GO TO (120,115),IDIR
      115 DO 116 J=1,NVARS
      DO 116 K=J,NVARS
      WECT(J,K) = W(J,K)
      CALL UPRFCT(NVARS,WECT,M)
      IF(M.NE.O) GO TO 125
      RETURN
      125 IFINE=1
      WRITE(6,900)
      FORMAT('1 THE WITHIN GROUPS MATRIX IS SINGULAR ')
      900 RETURN
      FUNCTION URAND(IRAND)
      THIS FUNCTION CALCULATES UNIFORMLY DISTRIBUTED RANDOM NUMBERS.
      BETWEEN 0 AND 1
      3**19 CONGRUENTIAL UNIFORM RANDCM NUMBER GENERATOR
      IRAND = IRAND*1162261467
      IF (IRAND.GT.O) GO TO 3
      IRAND = -IRAND
      URAND = FLOAT(IRAND)*0.4656612873E-9
      RETURN
      ENC
      SUBROUTINE UPRFCT(N,A,M)
      REPLACE UPPER TRIANGLE OF A SQUARE POSITIVE DEFINITE MATRIX A
      BY ITS CHOLESKI FACTOR
      DIMENSION A(20,20)
      CLEAR THE ERROR INDICATOR
      M=0
      N1=N-1

```

```

100 IF(N1) 230,100,100
    DO 220 K=1,N
    AKK=A(K,K)
    IF(K.EQ.1) GO TO 120
    DO 110 J=2,K
    110 AKK=AKK-A(J-1,K)**2
    120 IF(A(K,K)) 140,140,130
    IN THE CASE OF A COVARIANCE MATRIX AKK/A(K,K) IS 1-R**2 WHERE
    R IS MULT CORRELATION OF VARIABLE K WITH ALL PRECEDING VARIABLES.
    130 IF(AKK/A(K,K).GE..001) GO TO 150
    140 M=K
    150 AKK=0.
    AKK=SQRT(AKK)
    A(K,K)=AKK
    IF(K.EQ.N) GO TO 230
C
DO 220 I=K,N1
    AKI=A(K,I+1)
    IF(K.EQ.1) GO TO 190
    DO 180 J=2,K
    180 AKI=AKI-A(J-1,K)*A(J-1,I+1)
    190 IF(AKK) 210,200,210
    200 A(K,I+1)=0.0
    210 A(K,I+1)=AKI/AKK
    220 CONTINUE
    230 RETURN
    END
SUBROUTINE UTIRT(M,N,S,B)
C
C
    INVERSE OF UPPER (S) TRANSPOSED TIMES RECTANGLE B TRANSPOSED.
C
    DIMENSION S(20,20),B(1,20)
    DO 130 J=1,N
    DO 130 I=1,M
    SUM=0.0
    IF (S(I,I)) 90,130,90
    90 SUM = B(J,I)
    IM1 = I-1
    IF(IM1) 120,120,100
    DO 110 K=1,IM1
    110 SUM = SUM-S(K,I)*B(J,K)
    120 SUM = SUM/S(I,I)
    130 B(J,I) = SUM
    RETURN
    END
SUBROUTINE LTISUI (N,A,B)
C

```

```

C C      UPPER (B) TRANSPOSE INVERSE TIMES A TIMES UPPER (B) INVERSE.
C C      DIMENSION A(20,20),B(20,20)
C C      NOTE THAT POSTMULT IS CARRIED OUT ON FINAL VALUES LEFT BY PREMULT
C C      DC 200 I=1,N
C C      I1=I-1
C C      NOTE THAT PREMULT ON RIGHT HALF OF ROW IS SAME AS POSTMULT
C C      ON LOWER HALF OF COLUMN - EXCEPT FOR DIAG TERM WHICH IS THE
C C      FINAL VALUE LEFT BY PREMULT BEFORE POSTMULT
C C      DJ 130 J=I,N
C C      IF(I1) 120,120,100
C C      DO 110 K=1,I1
C C      100 A(J,I) = A(J,I) - B(K,I)*A(J,K)
C C      110 A(J,I) = A(J,I)/B(I,I)
C C      120 A(J,I) = A(J,I).EQ.0.0
C C      130 IF(B(I,I).EQ.0.0) A(J,I) = 0.0
C C      NOTE THAT ELEMENTS IN LEFT HALF OF ROW ARE FINAL FOR PREMULT
C C      NOTE THAT DIAG ELEMENT WAS PREVIOUSLY THE FINAL RESULT OF
C C      PREMULT. NOW WE MAKE THE FINAL RESULT OF POSTMULT.
C C      DC 200 J=1,I
C C      J1=J-1
C C      IF(J1) 160,160,140
C C      HORIZONTAL BRANCH OF INNER PRODUCT EXCLUDING DIAG TERM
C C      DO 150 K=1,J1
C C      140 A(I,J) = A(I,J) - B(K,I)*A(J,K)
C C      150 A(I,J) = A(I,J)
C C      160 IF(I1-J1) 190,170,170
C C      VERTICAL BRANCH OF INNER PRODUCT INCLUDING DIAG TERM
C C      DC 180 K=J1,I
C C      170 DC 180 K=J1,I
C C      180 A(I,J) = A(I,J) - B(K,I)*A(K,J)
C C      190 A(I,J) = A(I,J)/B(I,I)
C C      200 IF(B(I,I).EQ.0.0) A(I,J) = 0.0
C C      RETURN
C C      END
C C      SUBROUTINE EIGN(NN,A,EIG,VEC,INC)
C C      NN= SIZE OF MATRIX
C C      A= MATRIX (ONLY LOWER TRIANGLE IS USED + THIS IS DESTROYED)
C C      EIG = RETURNED EIGENVALUES IN ALGEBRAIC DESCENDING ORDER
C C      VEC = RETURNED EIGENVECTORS IN COLUMNS
C C      IND = ERROR RETURN INDICATOR
C C      0 FOR NORMAL RETURN
C C      1 SUM OF EIGENVALUES NOT EQUAL TO TRACE
C C      2 SUM OF EIGENVALUES SQUARED NOT EQUAL TO NORM
C C      3 BOTH OF THESE ERRORS
C C      DIMENSION A(20,20),GAMMA(20),BETA(20),EIG(20)
C C      DIMENSION W(20),VEC(20,20)
C C      THE FOLLOWING DIMENSIONED VARIABLES ARE EQUIVALENCED
C C      DIMENSION P(19),O(19)
C C      EQUIVALENCE (P(1),BETA(1)),(Q(1),BETA(1))
C C      DIMENSION IPOSV(20),IVPOS(20),IORD(20)
RRR10360
RRR10370
RRR10380
RRR10390
RRR10400
RRR10410
RRR10420
RRR10430
RRR10440
RRR10450
RRR10460
RRR10470
RRR10480
RRR10490
RRR10500
RRR10510
RRR10520
RRR10530
RRR10540
RRR10550
RRR10560
RRR10570
RRR10580
RRR10590
RRR10600
RRR10610
RRR10620
RRR10630
RRR10640
RRR10650
RRR10660
RRR10670
RRR10680
RRR10690
RRR10700
RRR10710
RRR10720
RRR10730
RRR10740
RRR10750
RRR10760
RRR10770
RRR10780
RRR10790
RRR10800
RRR10810
RRR10820
RRR10830

```



```

EQUIVALENCE (IPQSV(1), GAMMA(1)), (IVPOS(1), BETA(1)),
1(IORD(1), BETASQ(1))
N=NN
RESET ERROR RETURN INDICATOR
IND=0
IF(N.EQ. 0) GO TO 560
N1=N-1
N2=N-2
C COMPUTE THE TRACE AND EUCLIDIAN NCRM OF THE INPUT MATRIX
C LATER CHECK AGAINST SUM AND SUM OF SQUARES OF EIGENVALUES
ENCRM=0.
TRACE=0.
DO 110 J=1, N
DO 100 I=J, N
ENCRM=ENCRM+A(I, J)**2
TRACE=TRACE+A(I, J)
100 ENCRM=ENCRM-.5*A(I, J)**2
110 ENCRM=ENCRM+ENCRM
GAMMA(1)=A(1, 1)
IF(N2) 280, 270, 120
120 DO 260 NR=1, N2
E=A(NR+1, NR)
S=0.
DO 130 I=NR, N2
S=S+A(I+2, NR)**2
130 PREPARE FOR POSSIBLE BYPASS OF TRANSFORMATION
A(NR+1, NR)=0.
IF(S) 250, 250, 140
140 S=S+B*R
SGN=+1
IF(B) 150, 160, 160
150 SGNS=-1
SGRTS=SGRT(S)
D=SGN/(SGRTS+SGRTS)
TEMP=SGRT(.5+B*D)
W(NR)=TEMP
A(NR+1, NR)=TEMP
D=D/TEMP
B=-SGN*SGRTS
D IS FACTOR OF PROPORTIONALITY. NOW COMPUTE AND SAVE W VECTOR.
C EXTRA SINGLY SUBSCRIPTED W VECTOR USED FOR SPEED.
DO 170 I=NR, N2
TEMP=D*A(I+2, NR)
W(I+1)=TEMP
W(I+2, NR)=TEMP
17C PREMULTIPLY VECTOR W BY MATRIX A TO OBTAIN P VECTOR.
C SIMULTANEOUSLY ACCUMULATE DOT PRODUCT WP, (TYPE SCALAR K)
WTAW=0.

```

```

RRR10840
RRR10850
RRR10860
RRR10870
RRR10880
RRR10890
RRR10900
RRR10910
RRR10920
RRR10930
RRR10950
RRR10960
RRR10970
RRR10980
RRR10990
RRR11000
RRR11010
RRR11020
RRR11030
RRR11040
RRR11050
RRR11060
RRR11070
RRR11080
RRR11090
RRR11100
RRR11110
RRR11120
RRR11130
RRR11140
RRR11150
RRR11160
RRR11170
RRR11180
RRR11190
RRR11200
RRR11210
RRR11220
RRR11230
RRR11240
RRR11250
RRR11260
RRR11270
RRR11280
RRR11290
RRR11300
RRR11310

```

```

RRR111320
RRR111330
RRR111340
RRR111350
RRR111360
RRR111370
RRR111380
RRR111390
RRR111400
RRR111410
RRR111420
RRR111430
RRR111440
RRR111450
RRR111460
RRR111470
RRR111480
RRR111490
RRR111500
RRR111510
RRR111520
RRR111530
RRR111540
RRR111550
RRR111560
RRR111570
RRR111580
RRR111590
RRR111600
RRR111610
RRR111620
RRR111630
RRR111640
RRR111650
RRR111660
RRR111670
RRR111680
RRR111690
RRR111700
RRR111710
RRR111720
RRR111730
RRR111740
RRR111750
RRR111760
RRR111770
RRR111780
RRR111790

```

```

DC 220 I=NR,N1
SUM=0.
CC 180 J=NR,I
S1M=SUM+A(I+1,J+1)*W(J)
I1=I+1
IF(N1-I1) 210,190,190
DO 200 J=I,N1
S1M=SUM+A(J+1,I+1)*W(J)
P(I)=SUM
WTAW=WTAW+SUM*W(I)
P VECTOR AND SCALAR K NOW STORED. NEXT COMPLETE Q VECTOR
CC 230 I=NR,N1
Q(I)=P(I)-WTAW*W(I)
NOW FORM PAP MATRIX, REQUIRED PART
DO 240 J=NR,N1
QJ=Q(J)
WJ=W(J)
DO 240 I=J,N1
A(I+1,J+1)=A(I+1,J+1)-2.*(W(I)*QJ+WJ*Q(I))
240 A(I+1,J+1)=A(I+1,J+1)-2.*(W(I)*QJ+WJ*Q(I))
250 BETA(NR)=B
BETASQ(NR)=B*B
260 GAMMA(NR+1)=A(NR+1,NR+1)
B=A(N,N-1)
BETA(N-1)=B
BETASQ(N-1)=B*B
GAMMA(N)=A(N,N)
BETASQ(N)=0
ADJOIN AN IDENTITY MATRIX TO BE POSTMULTIPLIED BY ROTATIONS.
DO 300 I=1,N
DO 290 J=1,N
VEC(I,J)=0.
300 VEC(I,I)=1.
SUM=0.
NPAS=1
GO TO 400
S1M=SUM+SHIFT
CCSA=1.
G=GAMMA(1)-SHIFT
PP=G
PPES=PP*PP+BETASQ(1)
PPBR=SQRT(PPBS)
DO 370 J=1,M
CCSAP=CCSA
IF(PPBS .NE. 0.) GO TO 320
SINA=0.
SINA2=0.
CCSA=1.

```

```

RRR111800
RRR111810
RRR111820
RRR111830
RRR111840
RRR111850
RRR111860
RRR111870
RRR111880
RRR111890
RRR111900
RRR111910
RRR111920
RRR111930
RRR111940
RRR111950
RRR111960
RRR111970
RRR111980
RRR111990
RRR12000
RRR12010
RRR12020
RRR12030
RRR12040
RRR12050
RRR12060
RRR12070
RRR12080
RRR12090
RRR12100
RRR12110
RRR12120
RRR12130
RRR12140
RRR12150
RRR12160
RRR12170
RRR12180
RRR12190
RRR12200
RRR12210
RRR12220
RRR12230
RRR12240
RRR12250
RRR12260
RRR12270

```

```

32C GO TO 350
SINA=BETA(J)/PPBR
SINA2=BETASQ(J)/PPBS
COSA=PP/PPBR
POSTMULTIPLY IDENTITY BY P-TRANSPCSE MATRIX
NT=J+NPAS
IF(NT .GE. N) NT=N
DO 340 I=1,NT
TEMP=COSA*VEC(I,J)+SINA*VEC(I,J+1)
VEC(I,J+1)=-SINA*VEC(I,J)+COSA*VEC(I,J+1)
VEC(I,J)=TEMP
CIA=GAMMA(J+1)-SHIFT
U=SINA2*(G+DIA)
GAMMA(J)=G+U
G=DIA-U
PP=DIA*COSA-SINA*COSAP*BETA(J)
IF(J .NE. M) GO TO 360
BETA(J)=SINA*PP
BETASQ(J)=SINA2*PP*PP
GC TO 380
360 PPBS=PP*PP+BETASQ(J+1)
PPBR=SQRT(PPBS)
BETA(J)=SINA*PPBR
BETASQ(J)=SINA2*PPBS
37C GAMMA(M+1)=G
380 TEST FOR CONVERGENCE OF LAST DIAGONAL ELEMENT
NPAS=NPAS+1
IF(BETASQ(M) .GT. 1.E-21) GO TO 41C
EIG(M+1)=GAMMA(M+1)+SUM
BETA(M)=0.
BETASQ(M)=0.
M=M-1
IF(M .EQ. 0) GO TO 420
IF(BETASQ(M) .LE. 1.E-21) GO TO 390
TAKE ROOT OF CORNER 2 BY 2 NEAREST TO LOWER DIAGONAL IN VALUE
AS ESTIMATE OF EIGENVALUE TO USE FOR SHIFT
410 A2=GAMMA(M+1)
R2=.5*A2
R1=.5*GAMMA(M)
R12=R1+R2
DIF=R1-R2
TEMP=SQRT(DIF*DIF+BETASQ(M))
R1=R12+TEMP
R2=R12-TEMP
DIF=ABS(A2-R1)-ABS(A2-R2)
IF(DIF .LT. 0.) GO TO 420
SHIFT=R2
GO TO 310

```

```

420 SHIFT=R1
GO TO 310
43C EIG(I)=GAMMA(I)+SUM
C INITIALIZE AUXILIARY TABLES REQUIRED FOR REARRANGING THE VECTORS
DC 440 J=1,N
IPDSV(J)=J
IVPOS(J)=J
440 IORD(J)=J
C USE A TRANSPOSITION SORT TO ORDER THE EIGENVALUES
M=N
GO TO 470
450 DO 460 J=1,M
IF(EIG(J).GE. EIG(J+1)) GO TO 460
TEMP=EIG(J)
EIG(J)=EIG(J+1)
EIG(J+1)=TEMP
ITEMP=IORD(J)
IORD(J)=IORD(J+1)
IORD(J+1)=ITEMP
460 CONTINUE
470 M=M-1
IF(M.NE. 0) GO TO 450
IF(N1.EQ. 0) GO TO 500
DC 490 I=1,N1
NV=IORD(L)
NP=IPDSV(NV)
IF(NP.EQ. L) GO TO 490
LV=IVPOS(L)
IVPOS(NP)=LV
IPDSV(LV)=NP
DO 480 I=1,N
TEMP=VEC(I,L)
VEC(I,L)=VEC(I,NP)
VEC(I,NP)=TEMP
48C CONTINUE
49C ESUM=0.
500 ESSQ=0.
BACK TRANSFORM THE VECTORS OF THE TRIPLE DIAGONAL MATRIX
C DC 550 NRR=1,N
K=N1
51C K=K-1
IF(K.LE. 0) GO TO 540
SUM=0.
DO 520 I=K,N1
SLW=SUM+VEC(I+1,NRR)*A(I+1,K)
SUM=SUM+SLW
520 DO 530 I=K,N1
530 VEC(I+1,NRR)=VEC(I+1,NRR)-SUM*A(I+1,K)
RRR12280
RRR12290
RRR12300
RRR12310
RRR12320
RRR12330
RRR12340
RRR12350
RRR12360
RRR12370
RRR12380
RRR12390
RRR12400
RRR12410
RRR12420
RRR12430
RRR12440
RRR12450
RRR12460
RRR12470
RRR12480
RRR12490
RRR12500
RRR12510
RRR12520
RRR12530
RRR12540
RRR12550
RRR12560
RRR12570
RRR12580
RRR12590
RRR12600
RRR12610
RRR12620
RRR12630
RRR12640
RRR12650
RRR12660
RRR12670
RRR12680
RRR12690
RRR12700
RRR12710
RRR12720
RRR12730
RRR12740
RRR12750

```


RRR12840
RRR12830
RRR12810
RRR12800
RRR12790
RRR12780
RRR12770
RRR12760

```

//AIKSC179 JOB (1642,0526,RJ74),'AIKEN SMC 2132'
// EXEC FORTCLG,REGION.GO=180K
//FORT.SYSIN ON *

```

THIS METHOD OF COMPARING CLUSTERS OF OBJECTS BY MULTIPLE JUDGES
 WAS DESIGNED BY JOEL AIKEN WHO WAS GIVEN THE IDEA BY PROF J. HARTMAN.

ARRANGE DATA DECK SO THAT CARDS ARE IN THIS ORDER:

1. TITLE CARD FOLLOWED BY: (213)
2. # OF OBJECTS, # JUDGES (13)
3. # OF CLUSTERS BY FIRST JUDGE (13)
4. # OF OBJECTS IN EACH CLUSTER (1515)
5. (NOTE: 15 CLUSTERS IS MAX)
6. OBJECT ID NUMBERS FOR CLUSTERS (1615)
7. REPEAT AS NECESSARY, 2 CARDS PER CLUSTER.
8. RETURN TO STEP 2 FOR SUBSEQUENT JUDGES.

THE PROGRAM COMPUTES A SUM OF SQUARES VALUE AND DIVIDES THE SS
 BY THE BEST POSSIBLE SS AND THIS RATIO IS CALLED C.
 THE NUMBER AND SIZE OF CLUSTERS DETERMINE THE BEST SS.
 CSTAR IS COMPUTED AS (SS-LOW)/(BEST-LOW)

THE UPPER BOUND IS FOUND BY CONSIDERING THE # OF JUDGES,
 # OF CLUSTERS, AND # OF OBJECTS IN EACH CLUSTER AND FROM THAT
 DETERMINING THE BEST POSSIBLE MATCH MATRIX WHICH COULD EXIST.
 THIS PROCEDURE USES A MINIMAX APPROACH AND IS PERFORMED BY THE
 SUBROUTINE UPPER. THE NUMBER C IS THEN COMPUTED AS THE RATIO
 OF ACTUAL TO BEST SUM OF SQUARES. THIS PROCEDURE IS AN ATTEMPT TO
 ALLOW THE USER A STANDARD SCALE UPON WHICH TO COMPARE C VALUES
 BASED ON VARYING CLUSTER NUMBERS AND SIZES.
 SEE SUBROUTINE UPPER FOR DETAILS

```

COMMON IX1(100,100),IX2(100,15),ITEMP(100,100),NINCL(100),
1 NINC2(15),KCOUNT(100),NCLUS1,NCLUS2
DIMENSION IVY(50,10),L(10),TITLE(20),JROW(10)
SUM=0.0
READ(5,15) (TITLE(I),I=1,20)
FORMAT(20A4)
WRITE(6,16) (TITLE(I),I=1,20)
FORMAT(1,20A4,/)

```

15
 16
 CC

```

C      READ IN TOTAL # OBJECTS TO BE CLUSTERED AND
C      THE NUMBER OF JUDGES WHO WILL DO INDEPENDENT CLUSTERING
C
100    READ(5,100) NOBJ,NJUDGE
        FORMAT(2I3)
        DO 20 J=1,100
        DC 10 I=1,50
        ITEMP(I,J)=0
        CONTINUE
        CONTINUE

10      READ IN THE DATA FROM THE FIRST TWO JUDGES
20      READ(5,105) NCLUS1
        FORMAT(13)
        L(1)=NCLUS1
        READ(5,110) (NINC1(I),I=1,NCLUS1)
        FORMAT(15I5)
        WRITE(6,106) NOBJ,NJUDGE
        FORMAT(10 NUMBER OF OBJECTS TO BE CLUSTERED IS ',13,
110      11CX, NUMBER OF JUDGES IS ',13)
        WRITE(6,108)
        FORMAT(10,111,112,113,114,115,116,117,118,119,120,121,122,123,124,125,126,127,128,129,130,131,132,133,134,135,136,137,138,139,140,141,142,143,144,145,146,147,148,149,150,151,152,153,154,155,156,157,158,159,160,161,162,163,164,165,166,167,168,169,170,171,172,173,174,175,176,177,178,179,180,181,182,183,184,185,186,187,188,189,190,191,192,193,194,195,196,197,198,199,200,201,202,203,204,205,206,207,208,209,210,211,212,213,214,215,216,217,218,219,220,221,222,223,224,225,226,227,228,229,230,231,232,233,234,235,236,237,238,239,240,241,242,243,244,245,246,247,248,249,250,251,252,253,254,255,256,257,258,259,260,261,262,263,264,265,266,267,268,269,270,271,272,273,274,275,276,277,278,279,280,281,282,283,284,285,286,287,288,289,290,291,292,293,294,295,296,297,298,299,300,301,302,303,304,305,306,307,308,309,310,311,312,313,314,315,316,317,318,319,320,321,322,323,324,325,326,327,328,329,330,331,332,333,334,335,336,337,338,339,340,341,342,343,344,345,346,347,348,349,350,351,352,353,354,355,356,357,358,359,360,361,362,363,364,365,366,367,368,369,370,371,372,373,374,375,376,377,378,379,380,381,382,383,384,385,386,387,388,389,390,391,392,393,394,395,396,397,398,399,400,401,402,403,404,405,406,407,408,409,410,411,412,413,414,415,416,417,418,419,420,421,422,423,424,425,426,427,428,429,430,431,432,433,434,435,436,437,438,439,440,441,442,443,444,445,446,447,448,449,450,451,452,453,454,455,456,457,458,459,460,461,462,463,464,465,466,467,468,469,470,471,472,473,474,475,476,477,478,479,480,481,482,483,484,485,486,487,488,489,490,491,492,493,494,495,496,497,498,499,500,501,502,503,504,505,506,507,508,509,510,511,512,513,514,515,516,517,518,519,520,521,522,523,524,525,526,527,528,529,530,531,532,533,534,535,536,537,538,539,540,541,542,543,544,545,546,547,548,549,550,551,552,553,554,555,556,557,558,559,560,561,562,563,564,565,566,567,568,569,570,571,572,573,574,575,576,577,578,579,580,581,582,583,584,585,586,587,588,589,590,591,592,593,594,595,596,597,598,599,600,601,602,603,604,605,606,607,608,609,610,611,612,613,614,615,616,617,618,619,620,621,622,623,624,625,626,627,628,629,630,631,632,633,634,635,636,637,638,639,640,641,642,643,644,645,646,647,648,649,650,651,652,653,654,655,656,657,658,659,660,661,662,663,664,665,666,667,668,669,670,671,672,673,674,675,676,677,678,679,680,681,682,683,684,685,686,687,688,689,690,691,692,693,694,695,696,697,698,699,700,701,702,703,704,705,706,707,708,709,710,711,712,713,714,715,716,717,718,719,720,721,722,723,724,725,726,727,728,729,730,731,732,733,734,735,736,737,738,739,740,741,742,743,744,745,746,747,748,749,750,751,752,753,754,755,756,757,758,759,760,761,762,763,764,765,766,767,768,769,770,771,772,773,774,775,776,777,778,779,780,781,782,783,784,785,786,787,788,789,790,791,792,793,794,795,796,797,798,799,800,801,802,803,804,805,806,807,808,809,810,811,812,813,814,815,816,817,818,819,820,821,822,823,824,825,826,827,828,829,830,831,832,833,834,835,836,837,838,839,840,841,842,843,844,845,846,847,848,849,850,851,852,853,854,855,856,857,858,859,860,861,862,863,864,865,866,867,868,869,870,871,872,873,874,875,876,877,878,879,880,881,882,883,884,885,886,887,888,889,890,891,892,893,894,895,896,897,898,899,900,901,902,903,904,905,906,907,908,909,910,911,912,913,914,915,916,917,918,919,920,921,922,923,924,925,926,927,928,929,930,931,932,933,934,935,936,937,938,939,940,941,942,943,944,945,946,947,948,949,950,951,952,953,954,955,956,957,958,959,960,961,962,963,964,965,966,967,968,969,970,971,972,973,974,975,976,977,978,979,980,981,982,983,984,985,986,987,988,989,990,991,992,993,994,995,996,997,998,999,1000,1001,1002,1003,1004,1005,1006,1007,1008,1009,1010,1011,1012,1013,1014,1015,1016,1017,1018,1019,1020,1021,1022,1023,1024,1025,1026,1027,1028,1029,1030,1031,1032,1033,1034,1035,1036,1037,1038,1039,1040,1041,1042,1043,1044,1045,1046,1047,1048,1049,1050,1051,1052,1053,1054,1055,1056,1057,1058,1059,1060,1061,1062,1063,1064,1065,1066,1067,1068,1069,1070,1071,1072,1073,1074,1075,1076,1077,1078,1079,1080,1081,1082,1083,1084,1085,1086,1087,1088,1089,1090,1091,1092,1093,1094,1095,1096,1097,1098,1099,1100,1101,1102,1103,1104,1105,1106,1107,1108,1109,1110,1111,1112,1113,1114,1115,1116,1117,1118,1119,1120,1121,1122,1123,1124,1125,1126,1127,1128,1129,1130,1131,1132,1133,1134,1135,1136,1137,1138,1139,1140,1141,1142,1143,1144,1145,1146,1147,1148,1149,1150,1151,1152,1153,1154,1155,1156,1157,1158,1159,1160,1161,1162,1163,1164,1165,1166,1167,1168,1169,1170,1171,1172,1173,1174,1175,1176,1177,1178,1179,1180,1181,1182,1183,1184,1185,1186,1187,1188,1189,1190,1191,1192,1193,1194,1195,1196,1197,1198,1199,1200,1201,1202,1203,1204,1205,1206,1207,1208,1209,1210,1211,1212,1213,1214,1215,1216,1217,1218,1219,1220,1221,1222,1223,1224,1225,1226,1227,1228,1229,1230,1231,1232,1233,1234,1235,1236,1237,1238,1239,1240,1241,1242,1243,1244,1245,1246,1247,1248,1249,1250,1251,1252,1253,1254,1255,1256,1257,1258,1259,1260,1261,1262,1263,1264,1265,1266,1267,1268,1269,1270,1271,1272,1273,1274,1275,1276,1277,1278,1279,1280,1281,1282,1283,1284,1285,1286,1287,1288,1289,1290,1291,1292,1293,1294,1295,1296,1297,1298,1299,1300,1301,1302,1303,1304,1305,1306,1307,1308,1309,1310,1311,1312,1313,1314,1315,1316,1317,1318,1319,1320,1321,1322,1323,1324,1325,1326,1327,1328,1329,1330,1331,1332,1333,1334,1335,1336,1337,1338,1339,1340,1341,1342,1343,1344,1345,1346,1347,1348,1349,1350,1351,1352,1353,1354,1355,1356,1357,1358,1359,1360,1361,1362,1363,1364,1365,1366,1367,1368,1369,1370,1371,1372,1373,1374,1375,1376,1377,1378,1379,1380,1381,1382,1383,1384,1385,1386,1387,1388,1389,1390,1391,1392,1393,1394,1395,1396,1397,1398,1399,1400,1401,1402,1403,1404,1405,1406,1407,1408,1409,1410,1411,1412,1413,1414,1415,1416,1417,1418,1419,1420,1421,1422,1423,1424,1425,1426,1427,1428,1429,1430,1431,1432,1433,1434,1435,1436,1437,1438,1439,1440,1441,1442,1443,1444,1445,1446,1447,1448,1449,1450,1451,1452,1453,1454,1455,1456,1457,1458,1459,1460,1461,1462,1463,1464,1465,1466,1467,1468,1469,1470,1471,1472,1473,1474,1475,1476,1477,1478,1479,1480,1481,1482,1483,1484,1485,1486,1487,1488,1489,1490,1491,1492,1493,1494,1495,1496,1497,1498,1499,1500,1501,1502,1503,1504,1505,1506,1507,1508,1509,1510,1511,1512,1513,1514,1515,1516,1517,1518,1519,1520,1521,1522,1523,1524,1525,1526,1527,1528,1529,1530,1531,1532,1533,1534,1535,1536,1537,1538,1539,1540,1541,1542,1543,1544,1545,1546,1547,1548,1549,1550,1551,1552,1553,1554,1555,1556,1557,1558,1559,1560,1561,1562,1563,1564,1565,1566,1567,1568,1569,1570,1571,1572,1573,1574,1575,1576,1577,1578,1579,1580,1581,1582,1583,1584,1585,1586,1587,1588,1589,1590,1591,1592,1593,1594,1595,1596,1597,1598,1599,1600,1601,1602,1603,1604,1605,1606,1607,1608,1609,1610,1611,1612,1613,1614,1615,1616,1617,1618,1619,1620,1621,1622,1623,1624,1625,1626,1627,1628,1629,1630,1631,1632,1633,1634,1635,1636,1637,1638,1639,1640,1641,1642,1643,1644,1645,1646,1647,1648,1649,1650,1651,1652,1653,1654,1655,1656,1657,1658,1659,1660,1661,1662,1663,1664,1665,1666,1667,1668,1669,1670,1671,1672,1673,1674,1675,1676,1677,1678,1679,1680,1681,1682,1683,1684,1685,1686,1687,1688,1689,1690,1691,1692,1693,1694,1695,1696,1697,1698,1699,1700,1701,1702,1703,1704,1705,1706,1707,1708,1709,1710,1711,1712,1713,1714,1715,1716,1717,1718,1719,1720,1721,1722,1723,1724,1725,1726,1727,1728,1729,1730,1731,1732,1733,1734,1735,1736,1737,1738,1739,1740,1741,1742,1743,1744,1745,1746,1747,1748,1749,1750,1751,1752,1753,1754,1755,1756,1757,1758,1759,1760,1761,1762,1763,1764,1765,1766,1767,1768,1769,1770,1771,1772,1773,1774,1775,1776,1777,1778,1779,1780,1781,1782,1783,1784,1785,1786,1787,1788,1789,1790,1791,1792,1793,1794,1795,1796,1797,1798,1799,1800,1801,1802,1803,1804,1805,1806,1807,1808,1809,1810,1811,1812,1813,1814,1815,1816,1817,1818,1819,1820,1821,1822,1823,1824,1825,1826,1827,1828,1829,1830,1831,1832,1833,1834,1835,1836,1837,1838,1839,1840,1841,1842,1843,1844,1845,1846,1847,1848,1849,1850,1851,1852,1853,1854,1855,1856,1857,1858,1859,1860,1861,1862,1863,1864,1865,1866,1867,1868,1869,1870,1871,1872,1873,1874,1875,1876,1877,1878,1879,1880,1881,1882,1883,1884,1885,1886,1887,1888,1889,1890,1891,1892,1893,1894,1895,1896,1897,1898,1899,1900,1901,1902,1903,1904,1905,1906,1907,1908,1909,1910,1911,1912,1913,1914,1915,1916,1917,1918,1919,1920,1921,1922,1923,1924,1925,1926,1927,1928,1929,1930,1931,1932,1933,1934,1935,1936,1937,1938,1939,1940,1941,1942,1943,1944,1945,1946,1947,1948,1949,1950,1951,1952,1953,1954,1955,1956,1957,1958,1959,1960,1961,1962,1963,1964,1965,1966,1967,1968,1969,1970,1971,1972,1973,1974,1975,1976,1977,1978,1979,1980,1981,1982,1983,1984,1985,1986,1987,1988,1989,1990,1991,1992,1993,1994,1995,1996,1997,1998,1999,2000,2001,2002,2003,2004,2005,2006,2007,2008,2009,2010,2011,2012,2013,2014,2015,2016,2017,2018,2019,2020,2021,2022,2023,2024,2025,2026,2027,2028,2029,2030,2031,2032,2033,2034,2035,2036,2037,2038,2039,2040,2041,2042,2043,2044,2045,2046,2047,2048,2049,2050,2051,2052,2053,2054,2055,2056,2057,2058,2059,2060,2061,2062,2063,2064,2065,2066,2067,2068,2069,2070,2071,2072,2073,2074,2075,2076,2077,2078,2079,2080,2081,2082,2083,2084,2085,2086,2087,2088,2089,2090,2091,2092,2093,2094,2095,2096,2097,2098,2099,2100,2101,2102,2103,2104,2105,2106,2107,2108,2109,2110,2111,2112,2113,2114,2115,2116,2117,2118,2119,2120,2121,2122,2123,2124,2125,2126,2127,2128,2129,2130,2131,2132,2133,2134,2135,2136,2137,2138,2139,2140,2141,2142,2143,2144,2145,2146,2147,2148,2149,2150,2151,2152,2153,2154,2155,2156,2157,2158,2159,2160,2161,2162,2163,2164,2165,2166,2167,2168,2169,2170,2171,2172,2173,2174,2175,2176,2177,2178,2179,2180,2181,2182,2183,2184,2185,2186,2187,2188,2189,2190,2191,2192,2193,2194,2195,2196,2197,2198,2199,2200,2201,2202,2203,2204,2205,2206,2207,2208,2209,2210,2211,2212,2213,2214,2215,2216,2217,2218,2219,2220,2221,2222,2223,2224,2225,2226,2227,2228,2229,2230,2231,2232,2233,2234,2235,2236,2237,2238,2239,2240,2241,2242,2243,2244,2245,2246,2247,2248,2249,2250,2251,2252,2253,2254,2255,2256,2257,2258,2259,2260,2261,2262,2263,2264,2265,2266,2267,2268,2269,2270,2271,2272,2273,2274,2275,2276,2277,2278,2279,2280,2281,2282,2283,2284,2285,2286,2287,2288,2289,2290,2291,2292,2293,2294,2295,2296,2297,2298,2299,2300,2301,2302,2303,2304,2305,2306,2307,2308,2309,2310,2311,2312,2313,2314,2315,2316,2317,2318,2319,2320,2321,2322,2323,2324,2325,2326,2327,2328,2329,2330,2331,2332,2333,2334,2335,2336,2337,2338,2339,2340,2341,2342,2343,2344,2345,2346,2347,2348,2349,2350,2351,2352,2353,2354,2355,2356,2357,2358,2359,2360,2361,2362,2363,2364,2365,2366,2367,2368,2369,2370,2371,2372,2373,2374,2375,2376,2377,2378,2379,2380,2381,2382,2383,2384,2385,2386,2387,2388,2389,2390,2391,2392,2393,2394,2395,2396,2397,2398,2399,2400,2401,2402,2403,2404,2405,2406,2407,2408,2409,2410,2411,2412,2413,2414,2415,2416,2417,2418,2419,2420,2421,2422,2423,2424,2425,2426,2427,2428,2429,2430,2431,2432,2433,2434,2435,2436,2437,2438,2439,2440,2441,2442,2443,2444,2445,2446,2447,2448,2449,2450,2451,2452,2453,2454,2455,2456,2457,2458,2459,2460,2461,2462,2463,2464,2465,2466,2467,2468,2469,2470,2471,2472,2473,2474,2475,2476,2477,2478,2479,2480,2481,2482,2483,2484,2485,2486,2487,2488,2489,2490,2491,2492,2493,2494,2495,2496,2497,2498,2499,2500,2501,2502,2503,2504,2505,2506,2507,2508,2509,2510,2511,2512,2513,2514,2515,2516,2517,2518,2519,2520,2521,2522,2523,2524,2525,2526,2527,2528,2529,2530,2531,2532,2533,2534,2535,2536,2537,2538,2539,2540,2541,2542,2543,2544,2545,2546,2547,2548,2549,2550,2551,2552,2553,2554,2555,2556,2557,2558,2559,2560,2561,2562,2563,2564,2565,2566,2567,2568,2569,2570,2571,2572,2573,2574,2575,2576,2577,2578,2579,2580,2581,2582,2583,2584,2585,2586,2587,2588,2589,2590,2591,2592,2593,2594,2595,2596,2597,2598,2599,2600,2601,2602,2603,2604,2605,2606,2607,2608,2609,2610,2611,2612,2613,2614,2615,2616,2617,2618,2619,2620,2621,2622,2623,2624,2625,2626,2627,2628,2629,2630,2631,2632,2633,2634,2635,2636,2637,2638,2639,2640,2641,2642,2643,2644,2645,2646,2647,2648,2649,2650,2651,2652,2653,2654,2655,2656,2657,2658,2659,2660,2661,2662,2663,2664,2665,2666,2667,2668,2669,2670,2671,2672,2673,2674,2675,2676,2677,2678,2679,2680,2681,2682,2683,2684,2685,2686,2687,2688,2689,2690,2691,2692,2693,2694,2695,2696,2697,2698,2699,2700,2701,2702,2703,2704,2705,2706,2707,2708,27
```



```

3//, CSTAR = ',F10.4)
STOP
ENC

C
SUBROUTINE MERGE(IFLAG,ITCOL)
COMMON IX1(100,100),IX2(100,15),ITEMP(100,100),NINCL(100),
1NINC2(15),KOUNT1(100),NCLUS1,NCLUS2
IFLAG=0
ITCOL=1
KOUNT1(1)=1
BIG LOOP (4 DEEP)
C COMPARES EVERY POSSIBLE PAIR FOR A MATCH. ORDERED SO THAT COL 1 OF IX2
C IS CHECKED AGAINST ALL OF IX1(COL AT A TIME), THEN COL 2 OF IX1, ETC.
DO 500 N=1,NCLUS2
LIM4=NINCL(N)
DO 400 J=1,NCLUS1
LIM5=NINCL(J)
DO 200 I=1,LIM5
IF(IX2(M,N).NE.IX1(I,J)) GO TO 200
FOUND A MATCH: THAT WILL RECORD IT.
CALL TMAT(IFLAG,IX2(M,N),ITCOL,J,N)
GC TO 300
CONTINUE
CONTINUE
CONTINUE
ALTER WILL SET UP FOR NEXT JUDGE.
IF(N.EQ.NCLUS2) CALL ALTER(ITCOL)
CONTINUE
RETURN
ENC

200
300
400
C
500

C
SUBROUTINE TMAT(IFLAG,IVAL,ITCOL,J,N)
COMMON IX1(100,100),IX2(100,15),ITEMP(100,100),NINCL(100),
1NINC2(15),KOUNT1(100),NCLUS1,NCLUS2
CHECKS FOR FIRST TIME INTO HERE ON THIS MERGE.
IF (IFLAG.EQ.0) GO TO 150
CHECKS IF COLS ARE STILL THE SAME AS LAST MATCH: IF SO, WILL ADD ANOTHER
OBS TO THAT PARTICULAR MERGED CLUSTER.
IF(ICOL1.EQ.J.AND.ICOL2.EQ.N) GO TO 100
IF NEW COL, THEN NEW MERGED CLUSTER.
ITCOL=ITCOL+1
KOUNT1(1) IS # OF OBS IN I-TH MERGED CLUSTER.
KOUNT1(ITCOL)=1
GC TO 150
KOUNT1(ITCOL)=KOUNT1(ITCOL)+1
ITEMP(KOUNT1(ITCOL),ITCOL)=IVAL
SET UP FOR NEXT CALL. THESE WILL BE OLD COL #'S AGAINST WHICH WILL BE

```

```

C      CHECKED THE NEW COL #'S.
      ICOL1=J
      ICOL2=N
      IFLAG=1
      RETURN
      END

C      SUBROUTINE ALTER(ITCOL)
      COMMON IX1(100,100),IX2(100,15),ITEMP(100,100),NINCL(100),
      ININC(15),KOUNT1(100),NCLUS1,NCLUS2

      MAKES THE NAME OF THE MERGED SET: IX1.
      SO THE NEXT SET TO BE READ IN WILL BE A NEW IX2.

      NCLUS1=ITCOL
      DO 100 I=1,ITCOL
      NINCL(I)=KOUNT1(I)
      LIM6=KOUNT1(I)
      DO 50 J=1,LIM6
      IX1(J,I)=ITEMP(J,I)
      CONTINUE
      RETURN
      END

50
100
C      SUBROUTINE UPPER(IVY,L,JROW,NJUDGE,BEST,NOBJ)
      DIMENSION IVY(50,10),L(10),JROW(10)
      ISUM=0
      BEST=0.0

      THE MATRIX IVY IS CONSTRUCTED AS FOLLOWS:
      COLS = JUDGES
      ROWS = CLUSTERS
      THE (I,J)-TH ELEMENT OF THE MATRIX IS THE NUMBER OF OBJECTS IN
      THE I-TH CLUSTER OF THE J-TH JUDGE.

      FIND MAX IN EACH COLUMN

      DO 100 I=1,NJUDGE
      LIM1=L(I)
      DO 75 J=1,LIM1
      IF(J.NE.1) GO TO 60
      MAX=IVY(J,I)
      JROW(I)=J

```

```

60      IF (IVY(J,I),LE.MAX) GO TO 75
      MAX=IVY(J,I)
      JROW(I)=J
      CONTINUE
      CONTINUE

      NOW FIND MIN OF MAX'S AND SUBTRACT FROM EACH MAX

      DC 200 I=1,NJUDGE
      IF (I.EQ.1) MIN=IVY(JROW(1),I)
      IF (IVY(JROW(I),I).LT.MIN) MIN=IVY(JROW(I),I)
      CONTINUE

      ACCUMULATE SUM OF MINIMAX'S AND SUM OF SQUARES

      ISUM=ISUM+MIN
      BEST=BEST+MIN**2
      DO 300 I=1,NJUDGE
      IVY(JROW(I),I)=IVY(JROW(I),I)-MIN
      CONTINUE
      IF (ISUM.LT.NOBJ) GO TO 10
      RETURN
      END

//GO.SYSIN DD *

```

BIBLIOGRAPHY

1. Anderberg, M.R., Cluster Analysis For Applications, p. 3, Academic Press, 1973.
2. MacQueen, J., Some Methods for Classification and Analysis of Multivariate Observations, paper presented at Fifth Berkeley Symposium on Mathematical Statistics and Probability, University of California, 1965-1966.
3. McRae, D.J., Clustering Multivariate Observations, doctoral dissertation, University of North Carolina, Chapel Hill, N.C., 1973.
4. Friedman, H.P., and Rubin, J., "On Some Invariant Criteria for Grouping Data," Journal of the American Statistical Association, v. 62: 320, p. 1159-1178, December 1967.
5. Eisenheis, R.A. and Avery, R.B., Discriminant Analysis and Classification Procedures, p. 37, Lexington Books, 1972.
6. Anderson, T.W., An Introduction to Multivariate Statistical Analysis, Wiley, 1958.
7. Read, R.R., "A Study of SOF Data," unpublished report at Naval Postgraduate School, Monterey, Ca., 1 December 1977.
8. Lindsay, G.F., "On Constructing Interval Scales From Ordinal Judgments," July 1977; "Categorical Judgments: The Method of Successive Intervals", August 1976, unpublished reports at the Naval Postgraduate School; also Memorandum to Professor R.R. Read dated 2 October 1978, [NC4(551s)/pb], subject: "Results of OS 4207 Class Project in Scaling SOF Data."
9. Wilks, S.S., Mathematical Statistics, p. 573-601, Wiley, 1962.
10. Stanford University, Department of Statistics Technical Report No. 71, The Use of Faces to Represent Points in n-Dimensional Space Graphically, by Herman Chernoff, 27 December, 1971.
11. Chernoff, H., "The Use of Faces to Represent Points in k-Dimensional Space Graphically," Journal of the American Statistical Association, V. 68: 342, p. 361-368, June, 1973.

12. Lake, G.E., A Graphical Representation of Multivariate Data: Soviet Foreign Policy in Sub-Saharan Africa, Master's Thesis, Naval Postgraduate School, Monterey, Ca., 1977.
13. Chernoff, H. and Rizvi, M.H., "Error of Random Permutations of Features in Representing Multivariate Data by Faces," Journal of the American Statistical Association, V. 70:351, p. 548-554, September 1975.
14. Siegel, S., Nonparametric Statistics, p. 201, McGraw-Hill, 1956.

INITIAL DISTRIBUTION LIST

	No. Copies
1. Defense Documentation Center Cameron Station Alexandria, Virginia 22314	2
2. Library, Code 0142 Naval Postgraduate School Monterey, California 93940	2
3. Department Chairman, Code 55 Department of Operations Research Naval Postgraduate School Monterey, California 93940	1
4. Professor R.R. Read, Code 55Re Department of Operations Research Naval Postgraduate School Monterey, California 93940	1
5. Professor J.R. Borsting, Code 01 Provost Naval Postgraduate School Monterey, California 93940	1
6. Professor D.E. Kirk, Code 62Ki Chairman Department of Electrical Engineering Naval Postgraduate School Monterey, California 93940	1
7. LCDR Joel W. Aiken 555 "C" Avenue Coronado, California 92118	1
8. Scholarship Committee Naval Postgraduate School Monterey, California 93940	1
9. Professor G.F. Lindsay, Code 55Ls Department of Operations Research Naval Postgraduate School Monterey, California 93940	1
10. D.J. McRae Cal Test Bureau McGraw-Hill Monterey, California 93940	1